

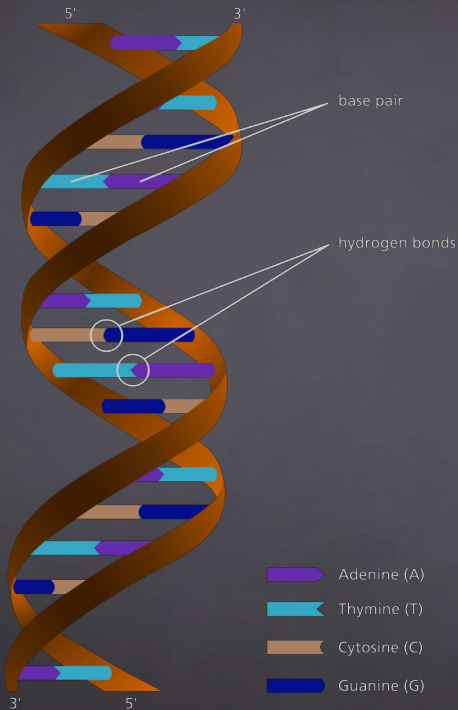
# Lecture: Graph-Based Genome Scaffolding

Mathias Weller

`mathias.weller@univ-mlv.fr`

Montpellier, 2017



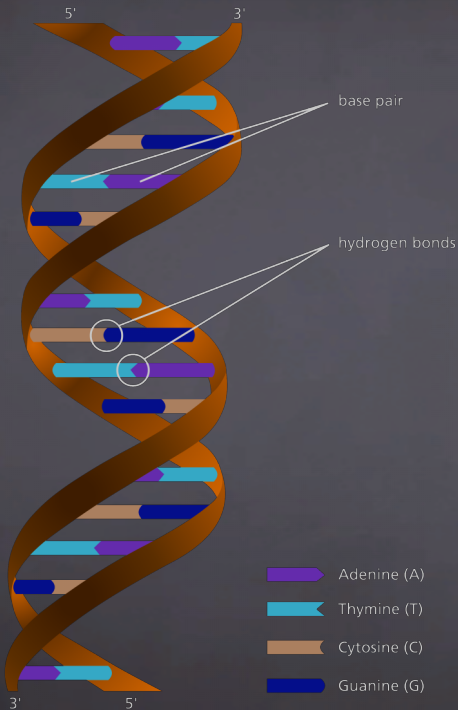


## DNA

- double strand
- inside nucleus (safe)

## RNA

- single strand
- outside nucleus
- transfers genetic code
- Thymine (T) → Uracil (U)



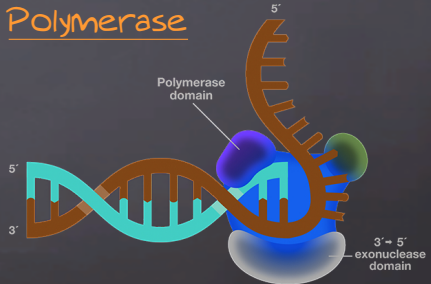
## DNA

- double strand
- inside nucleus (safe)

## RNA

- single strand
- outside nucleus
- transfers genetic code
- Thymine (T) → Uracil (U)

## Polymerase



# Sanger Sequencing

[Sanger et al '77]

CCTGGACGGGTCAGACATGACAGTGGCCCCAAGATTCACAAGATCGTATCTCAATACAGTAAACGAGCAAT  
GGACCTGCCCAGTCTGTACTGTCACCGGGGTTCTAAGTGTTCTAGCATAGAGTTATGTCATTTGCTCGTTA

# Sanger Sequencing

[Sanger et al '77]

CCTGGACGGGTCAGACATGACAGTGGCCCCAAGATTCACAAGATCGTATCTCAATACAGTAAACGAGCAAT  
GGACCTGCCCAGTCTGTACTGTCACCGGGGTTCTAAGTGTTCCTAGCATAGAGTTATGTCATTTGCTCGTTA

## Sanger Sequencing

1. split helix & create thousands of copies

# Sanger Sequencing

[Sanger et al '77]

CCTGGACGGGTCAGACATGACAGTGGCCCCAAGATTCACAAGATCGTATCTCAATACAGTAAACGAGCAAT

## Sanger Sequencing

1. split helix & create thousands of copies

# Sanger Sequencing

[Sanger et al '77]

CCTGGACGGGTCAGACATGACAGTGGCCCCAAGATTCACAAGATCGTATCTCAATACAGTAAACGAGCAAT

## Sanger Sequencing

1. split helix & create thousands of copies
2. add polymerase & floating Bases: A C G T
3. add a special Base: A\* (polymerase cannot extend)

# Sanger Sequencing

[Sanger et al '77]

CCTGGACGGGTCAGACATGACAGTGGCCCCAAGATTCACAAGATCGTATCTCAATACAGTAAACGAGCAAT

## Sanger Sequencing

1. split helix & create thousands of copies
2. add polymerase & floating Bases: A C G T
3. add a special Base: A\* (polymerase cannot extend)
4. stir & let polymerase act

# Sanger Sequencing

[Sanger et al '77]

CCTGGACGGGTCAGACATGACAGTGGCCCCAAGATTCACAAGATCGTATCTCAATACAGTAAACGAGCAAT  
GGA\*

## Sanger Sequencing

1. split helix & create thousands of copies
2. add polymerase & floating Bases: A C G T
3. add a special Base: A\* (polymerase cannot extend)
4. stir & let polymerase act

# Sanger Sequencing

[Sanger et al '77]

CCTGGACGGGTCAGACATGACAGTGGCCCCAAGATTCACAAGATCGTATCTCAATACAGTAAACGAGCAAT  
GGA\*  
GGACCTGCCCA\*

## Sanger Sequencing

1. split helix & create thousands of copies
2. add polymerase & floating Bases: A C G T
3. add a special Base: A\* (polymerase cannot extend)
4. stir & let polymerase act

# Sanger Sequencing

[Sanger et al '77]

CCTGGACGGGTCAGACATGACAGTGGCCCCAAGATTCACAAGATCGTATCTCAATACAGTAAACGAGCAAT  
GGA\*  
GGACCTGCCCA\*  
GGACCTGCCCAGTCTGTA\*

## Sanger Sequencing

1. split helix & create thousands of copies
2. add polymerase & floating Bases: A C G T
3. add a special Base: A\* (polymerase cannot extend)
4. stir & let polymerase act

# Sanger Sequencing

[Sanger et al '77]

CCTGGACGGGTCAGACATGACAGTGGCCCCAAGATTACAAAGATCGTATCTCAATACAGTAAACGAGCAAT  
GGA\*  
GGACCTGCCCCA\*  
GGACCTGCCCAGTCTGTA\*

## Sanger Sequencing

1. split helix & create thousands of copies
2. add polymerase & floating Bases: A C G T
3. add a special Base: A\* (polymerase cannot extend)
4. stir & let polymerase act
5. measure the length of each fragment  
→ each length is the position of a T in the template

# Sanger Sequencing

[Sanger et al '77]

CCTGGACGGGTCAGACATGACAGTGGCCCCAAGATTCACAAGATCGTATCTCAATACAGTAAACGAGCAAT  
GGA\*  
GGACCTGCCCA\*  
GGACCTGCCCAGTCTGTA\*

## Sanger Sequencing

1. split helix & create thousands of copies
2. add polymerase & floating Bases: A C G T
3. add a special Base: A\* (polymerase cannot extend)
4. stir & let polymerase act
5. measure the length of each fragment  
~> each length is the position of a T in the template

## Problem

unreliable after a couple hundred bp  
~> chop up DNA into pieces and read those

# Next Generation Sequencing (illumina)



# Next Generation Sequencing (illumina)

ACTCA.....ACCTC

I. chop DNA into smaller pieces

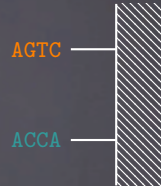
# Next Generation Sequencing (illumina)

TGGTACTCA.....ACCTCTCAG

1. chop DNA into smaller pieces
2. add anchors to each end of each piece

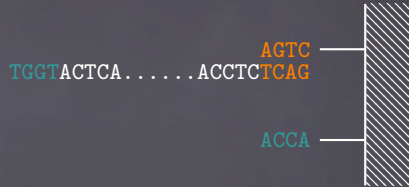
# Next Generation Sequencing (illumina)

TGGTACTCA.....ACCTCTCAG



1. chop DNA into smaller pieces
2. add anchors to each end of each piece
3. "flow cell" containing anchor places

# Next Generation Sequencing (illumina)



1. chop DNA into smaller pieces
2. add anchors to each end of each piece
3. "flow cell" containing anchor places
4. strand anchors its two ends to two anchor places

# Next Generation Sequencing (illumina)



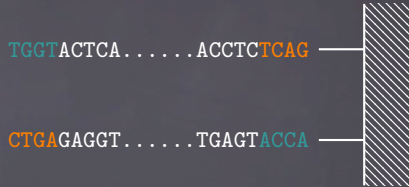
1. chop DNA into smaller pieces
2. add anchors to each end of each piece
3. "flow cell" containing anchor places
4. strand anchors its two ends to two anchor places

# Next Generation Sequencing (illumina)



1. chop DNA into smaller pieces
2. add anchors to each end of each piece
3. "flow cell" containing anchor places
4. strand anchors its two ends to two anchor places
5. polymerase completes the strand into double-strand

# Next Generation Sequencing (illumina)

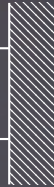


1. chop DNA into smaller pieces
2. add anchors to each end of each piece
3. "flow cell" containing anchor places
4. strand anchors its two ends to two anchor places
5. polymerase completes the strand into double-strand
6. double strand is denaturized into single strands

# Next Generation Sequencing (illumina)

TGGTACTCA . . . . . ACCTCTCAG

CTGAGAGGT . . . . . TGAGTACCA

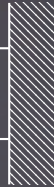


1. chop DNA into smaller pieces
2. add anchors to each end of each piece
3. "flow cell" containing anchor places
4. strand anchors its two ends to two anchor places
5. polymerase completes the strand into double-strand
6. double strand is denaturized into single strands
7. rinse, repeat (last 3 steps) until flow chip is "full"

# Next Generation Sequencing (illumina)

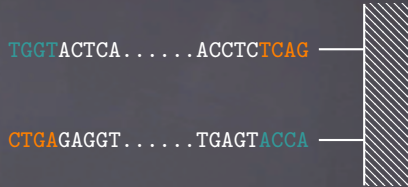
TGGTACTCA . . . . . ACCTCTCAG

CTGAGAGGT . . . . . TGAGTACCA



1. chop DNA into smaller pieces
2. add anchors to each end of each piece
3. "flow cell" containing anchor places
4. strand anchors its two ends to two anchor places
5. polymerase completes the strand into double-strand
6. double strand is denaturized into single strands
7. rinse, repeat (last 3 steps) until flow chip is "full"
8. read all strands from their anchor points outwards

# Next Generation Sequencing (illumina)



1. chop DNA into smaller pieces
  2. add anchors to each end of each piece
  3. "flow cell" containing anchor places
  4. strand anchors its two ends to two anchor places
  5. polymerase completes the strand into double-strand
  6. double strand is denaturized into single strands
  7. rinse, repeat (last 3 steps) until flow chip is "full"
  8. read all strands from their anchor points outwards
- ~> Paired-End reads (distance between reads = "insert size")

# Sequence Assembly: Overview

GGCCCTGAACTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

# Sequence Assembly: Overview

GGCCCTGAACCTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTGCGGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GGCCCTGAACCTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTGCGGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
ACTTCGC CGACACTCCTTGGGTTTT CTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

# Sequence Assembly: Overview

GCCCCCTGAACTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTGCGGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACTT CGACACTCCTTGGGTTTT CTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
ACTTCGC GGTTCCT GGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACTTCGC CGACACTCCTTGGGTTTT GGTTCCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

# Sequence Assembly: Overview

GCCCCTGAACTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTGCGGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCCTGAACTT ACTTCGC CGACACTCCTTGGGTTTT CTAGGCCATTGATTGCGGGTCCAGGTGCTGTGTCAACGAC  
GCCCCTGAACTTCGC CGACACTCCTTGGGTTTT GGTTCCTCT GGTCCAGGTGCTGTGTCAACGAC  
GTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

Problem 1: parts of the sequence might not be covered by reads

# Sequence Assembly: Overview

GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACCTT CGACACTCCTTGGGTTTT CTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
ACTTCGC GGTCTCT GGTCCAGGTGCTGTCAACGAC  
TCGCTAGGGTTCTCTAACGA TTTACGTCGCGG CGAC

Goal: reconstruct sequence

Idea: overlap reads

Problem 1: parts of the sequence might not be covered by reads  
~> sequence with "high coverage"

# Sequence Assembly: Overview

GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACCTT CGACACTCCTTGGGTTTT CTAGGCCATTGATTGCGGGTC  
ACTTCGC GGTTCCT GGTCCAGGTGCTGTCAACGAC  
TCGCTAGGGTTCTCTAACGA TTTACGTCGCGG CGAC  
GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

Problem 1: parts of the sequence might not be covered by reads  
~> sequence with "high coverage"

# Sequence Assembly: Overview

GCCCCCTGAACTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTGCGGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACTT CGACACTCCTTGGGTTTT CTAGGCCATTGATTGCGGGTCC  
ACTTCGC GGTTCCTCT GGTCCAGGTGCTGTCAACGAC  
TCGCTAGGGTTCTCTAACGA TTTACGTGCGGG CGAC  
GCCCCCTGAACTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTGCGGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

Problem 2: Shortest Common Superstring is NP-hard

~> "Overlap-Layout-Consensus" assemblers

# Sequence Assembly: Overview

GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACCTT CGACACTCCTTGGGTTTT CTAGGCCATTGATTGCGGGTC  
ACTTCGC GGTTCCTT GGTCCAGGTGCTGTCAACGAC  
TCGCTAGGGTTCTCTAACGA TTACGTCGCGG CGAC  
GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

Problem 2: Shortest Common Superstring is NP-hard

~> "Overlap-Layout-Consensus" assemblers

1. produce best pairwise overlaps
2. layout the reads according to the overlaps
3. for each position, compute consensus base

# Sequence Assembly: Overview

GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACCTC CGACACTCCTTGGGTTTT CTAGGCCATTGATTGCGGGTC  
ACTTCGC GGTTCCTCT GGTCCAGGTGCTGTCAACGAC  
TCGCTAGGGTTCTCTAACGA TTTACGTCGCGG CGAC  
GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

Problem 2: Shortest Common Superstring is NP-hard

~> "Overlap-Layout-Consensus" assemblers

1. produce best pairwise overlaps
2. layout the reads according to the overlaps
3. for each position, compute consensus base

# Sequence Assembly: Overview

GCCCCCTGAACTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACTT CGACACTCCTTGGGTTTT CTAGGCCATTGATTGCGGGTC  
ACTTCGC GGTTCCTT GGTCCAGGTGCTGTCAACGAC  
TCGCTAGGGTTCTCTAACGA TTTACGTCGCGG CGAC  
GCCCCCTGAACTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

Problem 2: Shortest Common Superstring is NP-hard

~> "Overlap-Layout-Consensus" assemblers

Problem:  $\Theta(n^2)$  too slow in practice

~> DeBruijn-graph Based assembly

# Sequence Assembly: Overview

GCCCTGAACCTTCGCTAGCTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCTGAACCTT CGACACTCCTTGGGTTTT CTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
ACTTCGC GGTTCCT GGTCCAGGTGCTGTCAACGAC  
TCGCTAGCTTCTCTAACGA TTTACGTCGCGG CGAC  
GCCCTGAACCTTCGCTAGCTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

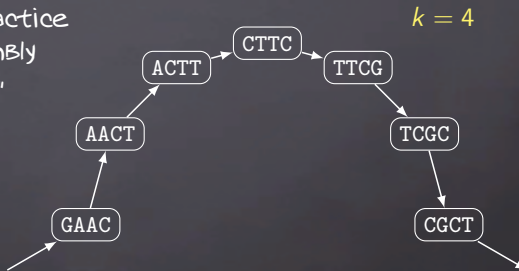
Problem 2: Shortest Common Superstring is NP-hard

~> "Overlap-Layout-Consensus" assemblers

Problem:  $\Theta(n^2)$  too slow in practice

~> DeBruijn-graph Based assembly

1. chop all reads into " $k$ -mers"
2. builds overlap graph ("DeBruijn graph")
3. find path using all overlaps



# Sequence Assembly: Overview

GGCCCTGAACCTTCGCTAGCTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GGCCCTGAACCTT CGACACTCCTTGGGTTTT CTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
ACTTCGC GGTTCCTT GGTCCAGGTGCTGTCAACGAC  
TCGCTAGCTTCTCTAACGA TTTACGTCGCGG CGAC  
GGCCCTGAACCTTCGCTAGCTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

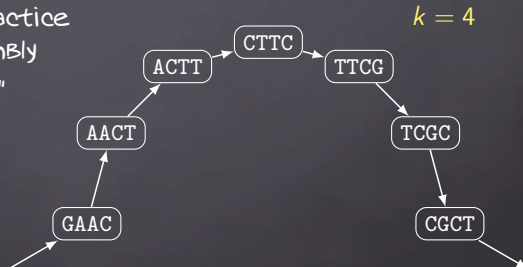
Problem 2: Shortest Common Superstring is NP-hard

~> "Overlap-Layout-Consensus" assemblers

Problem:  $\Theta(n^2)$  too slow in practice

~> DeBruijn-graph Based assembly

1. chop all reads into " $k$ -mers"
2. builds overlap graph ("DeBruijn graph")
3. find Eulerian path



# Sequence Assembly: Overview

GCGGCTGAACTTCGCTAGCTTCTCTAACGACACTCCTTGGGTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCGGCTGAACTTCTGACACTCCTTGGGTTTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
ACTTCGC  
TCGCTAGCTTCTCTAACGA  
TTTACGTCGCGG  
GGTTCTCT  
GGTCCAGGTGCTGTCAACGAC  
CGAC

Goal: reconstruct sequence

Idea: overlap reads

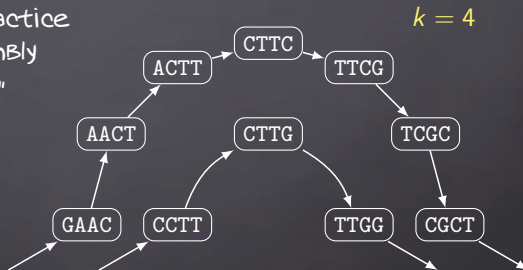
Problem 2: Shortest Common Superstring is NP-hard

→ "Overlap-Layout-Consensus" assemblers

Problem:  $\Theta(n^2)$  too slow in practice

→ DeBruijn-graph Based assembly

1. chop all reads into " $k$ -mers"
2. builds overlap graph ("DeBruijn graph")
3. find Eulerian path



# Sequence Assembly: Overview

GCCCTGAACCTTCGCTAGCTTCTCTAACGACACTCCTTGGGTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCTGAACCTT CGACACTCCTTGGGTTT CTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
ACTTCGC GGTTCCTT GGTCCAGGTGCTGTCAACGAC  
TCGCTAGCTTCTCTAACGA TTTACGTCGCGG CGAC  
GCCCTGAACCTTCGCTAGCTTCTCTAACGACACTCCTTGGGTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

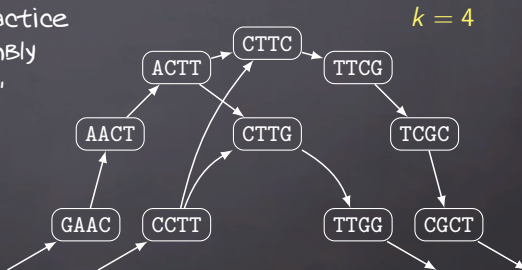
Problem 2: Shortest Common Superstring is NP-hard

→ "Overlap-Layout-Consensus" assemblers

Problem:  $\Theta(n^2)$  too slow in practice

→ DeBruijn-graph Based assembly

1. chop all reads into " $k$ -mers"
2. builds overlap graph ("DeBruijn graph")
3. find Eulerian path



# Sequence Assembly: Overview

GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACCTT CGACACTCCTTGGGTTTT CTAGGCCATTGATTGCGGGTC  
ACTTCGC GGTTCCTT GGTCCAGGTGCTGTCAACGAC  
TCGCTAGGGTTCTCTAACGA TTTACGTCGCGG CGAC  
GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

Problem 3: repeats (common in DNA) make assembly ambiguous

# Sequence Assembly: Overview

GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
ACTTCGC CGACACTCCTTGGGTTTTT GGTTCCTCT GGTCCAGGTGCTGTCAACGAC  
TCGCTAGGGTTCTCTAACGA TTTACGTCGCGG CGAC  
GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

Problem 3: repeats (common in DNA) make assembly ambiguous

# Sequence Assembly: Overview

GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
ACTTCGC GGTTCCTCT TCGCTAGGGTTCTCTAACGA TTTACGTCGCGG GGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGG CTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

Problem 3: repeats (common in DNA) make assembly ambiguous

# Sequence Assembly: Overview

GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
ACTTCGC GGTTCTCT CGACACTCCTTGGGTTTTTACGTCGCGGGTTCTTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
TCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

Problem 3: repeats (common in DNA) make assembly ambiguous  
~> end product is a set of "contiguous regions"

# Sequence Assembly: Overview

GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
ACTTCGC GGTTCTCT CGACACTCCTTGGGTTTTTACGTCGCGGGTTCTTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
TCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

Problem 3: repeats (common in DNA) make assembly ambiguous  
~> end product is a set of "contiguous regions"

Problem: "contig soup" not very useful

# Sequence Assembly: Overview

GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
ACTTCGC GGTTCCTCT TCGCTAGGGTTCTCTAACGA TTTACGTCGCGG GGTCCAGGTGCTGTCAACGAC  
GCCCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGG CTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

Problem 3: repeats (common in DNA) make assembly ambiguous  
~> end product is a set of "contiguous regions"

Problem: "contig soup" not very useful

But: we have paired-end information!

# Sequence Assembly: Overview

GGCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
GGCCCTGAACCTTCGCTAGGGTTCTCTAACGACACTCCTTGGGTTTTTACGTCGCGGGTTCTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
ACTTCGC GGTTCCTCT TCGCTAGGGTTCTCTAACGA TTTACGTCGCGG NNNNCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC  
CGACACTCCTTGGGTTTTTCTAGGCCATTGATTGCGGGTCCAGGTGCTGTCAACGAC

Goal: reconstruct sequence

Idea: overlap reads

Problem 3: repeats (common in DNA) make assembly ambiguous  
~> end product is a set of "contiguous regions"

Problem: "contig soup" not very useful

But: we have paired-end information!

# Genome Scaffolding: Previous Work

Goal: order & orient contigs

Idea: use pairing information on reads to "link" contigs together

## - SOPRA

[Dayarian, Michael, Sengupta, BMC Bioinf. 11, '10]

- ▶ removes reads in high-coverage area (likely repeats)
- ▶ orientation step (heuristic) + ordering step (heuristic)
- ▶ coded in **Pearl** (!!!)
- ▶ (observed **sparse** contig graph)

# Genome Scaffolding: Previous Work

Goal: order & orient contigs

Idea: use pairing information on reads to "link" contigs together

- SOPRA

[Dayarian, Michael, Sengupta, BMC Bioinf. 11, '10]

- SSPACE

[Boetzer & al., Bioinf. 27(4), '11]

- ▶ heuristic contig extension
- ▶ "reasonable time"

# Genome Scaffolding: Previous Work

Goal: order & orient contigs

Idea: use pairing information on reads to "link" contigs together

- SOPRA
- SSPACE
- OPERA

[Dayarian, Michael, Sengupta, BMC Bioinf. 11, '10]

[Boetzer & al., Bioinf. 27(4), '11]

[Gao, Sung, Nagaraja, JCB. 18(11), '11]

- ▶  $n^{p+O(1)}$  time ( $p = \# \text{edge-deletions } (\geq \text{feedback edge set})$ )
- ▶ most work done by a heuristic "graph contraction"

# Genome Scaffolding: Previous Work

Goal: order & orient contigs

Idea: use pairing information on reads to "link" contigs together

- SOPRA
- SSPACE
- OPERA
- GRASS

[Dayarian, Michael, Sengupta, BMC Bioinf. 11, '10]

[Boetzer & al., Bioinf. 27(4), '11]

[Gao, Sung, Nagaraja, JCB. 18(11), '11]

[Gritsenko & al., Bioinf. 28(11), '12]

- ▶ Mixed-Integer Quadratic Programming
- ▶ deals with uncertain data (slack variables)
  - ~> "intractable even for small # of contigs"
- ▶ heuristic workaround:
  - ▶ solve relaxed formulation & use slack values ~> ILP

# Genome Scaffolding: Previous Work

Goal: order & orient contigs

Idea: use pairing information on reads to "link" contigs together

- SOPRA
- SSPACE
- OPERA
- GRASS
- SCARPA

[Dayarian, Michael, Sengupta, BMC Bioinf. 11, '10]

[Boetzer & al., Bioinf. 27(4), '11]

[Gao, Sung, Nagaraja, JCB. 18(11), '11]

[Gritsenko & al., Bioinf. 28(11), '12]

[Donmez, Brudno, Bioinf. 29(4), '13]

- ▶ orientation step: use FPT algo for **Odd Cycle Transversal**
- ▶ ordering step: heuristic

# Genome Scaffolding: Previous Work

Goal: order & orient contigs

Idea: use pairing information on reads to "link" contigs together

- SOPRA
- SSPACE
- OPERA
- GRASS
- SCARPA
- ...

[Dayarian, Michael, Sengupta, BMC Bioinf. 11, '10]

[Boetzer & al., Bioinf. 27(4), '11]

[Gao, Sung, Nagaraja, JCB. 18(11), '11]

[Gritsenko & al., Bioinf. 28(11), '12]

[Donmez, Brudno, Bioinf. 29(4), '13]

[Huson & al., JACM, '02][Nieuwerburgh & al., NAR, '12]

# Graph-Based Scaffolding

AACGACACTCCTTGGGTTTTACGTGCGG

GTTAATGTCCGAGCATAAACTCTGGTTGGC

GTACTGAACTTGGGTTCCATAGGAACCAGA

CTAGGCCAATTGATTGCCGGTCCAGGTGCTG

AGAGCTTGACAGTAACACATTTAGGAGCACGCG

# Graph-Based Scaffolding

AACGACACTCCTGGGTTTTACGTCGCGG

GTTAATGTCCGAGCATAAACTCTGGTTGGC

GTACTGAAC TTGGGTTCCATAGGAACCAGA

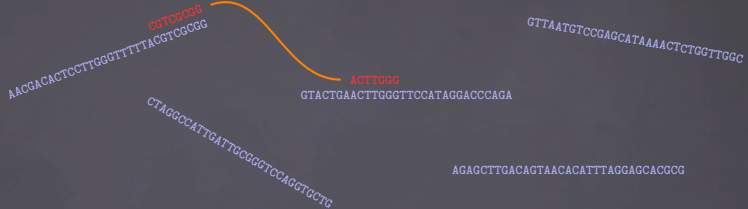
CTAGGCCAATTGATTGCGGGTCCAGGTGCTG

AGAGCTTGACAGTAACACATTTAGGAGCACGCG

## Strategy

- I. map reads into contigs

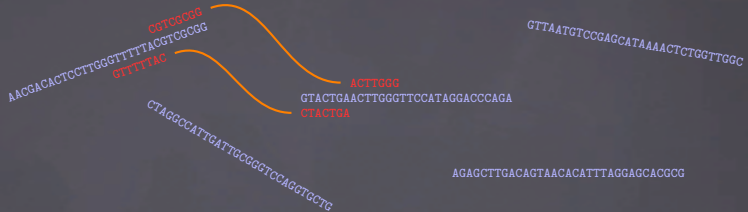
# Graph-Based Scaffolding



## Strategy

I. map reads into contigs

# Graph-Based Scaffolding



## Strategy

- I. map reads into contigs

# Graph-Based Scaffolding

AACGACACTCCTTGGGTTTTACGTGCGG

GTTAATGTCCGAGCATAAACTCTGGTTGGC

GTACTGAACTTGGGTTCCATAGGAACCAGA

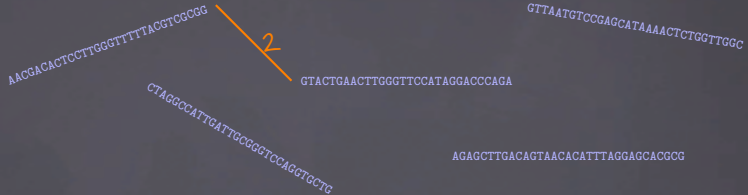
CTAGGCCAATTGATTGCGGCTCCAGGTGCTG

AGAGCTTGACAGTAACACATTTAGGAGCACGCG

## Strategy

1. map reads into contigs
2. pair contigs according to read-pairing (weighted)

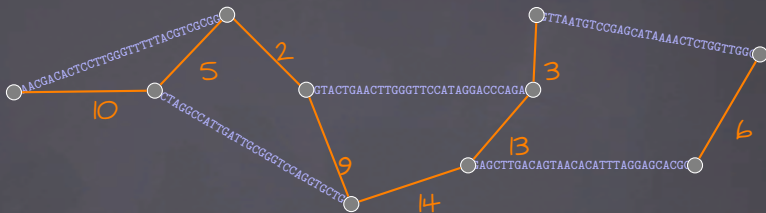
# Graph-Based Scaffolding



## Strategy

1. map reads into contigs
2. pair contigs according to read-pairing (weighted)

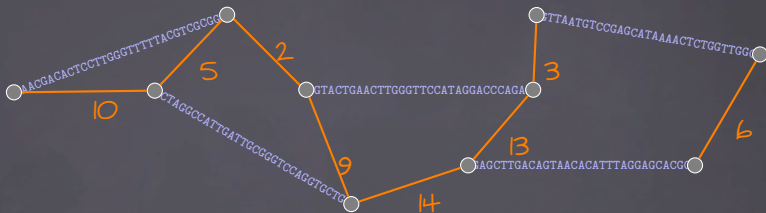
# Graph-Based Scaffolding



## Strategy

1. map reads into contigs
2. pair contigs according to read-pairing (weighted)

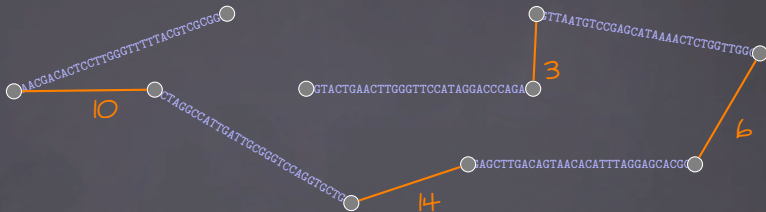
# Graph-Based Scaffolding



## Strategy

1. map reads into contigs
2. pair contigs according to read-pairing (weighted)
3. cover "scaffold graph" with (heavy) alternating paths  
each path corresponds to a chromosome

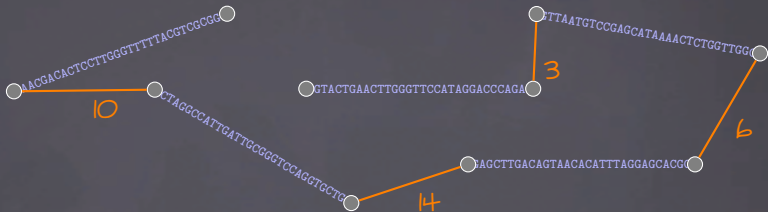
# Graph-Based Scaffolding



## Strategy

1. map reads into contigs
2. pair contigs according to read-pairing (weighted)
3. cover "scaffold graph" with (heavy) alternating paths  
each path corresponds to a chromosome

# Graph-Based Scaffolding



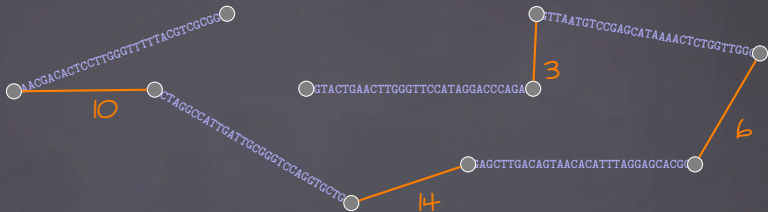
## Scaffolding

Input: Graph  $G$ , perfect matching  $M$ , weights  $w, k, \sigma_p \in \mathbb{N}$

Question: Can  $G$  be covered by  
 $\leq \sigma_p$  alternating paths

of total weight  $\geq k$ ?

# Graph-Based Scaffolding

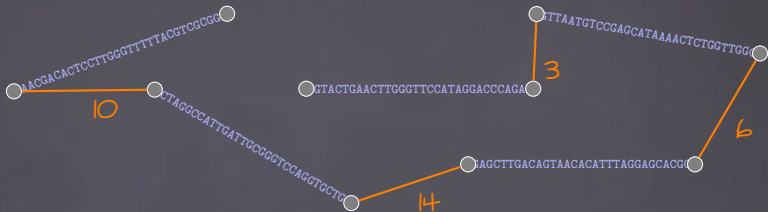


## Scaffolding

Input: Graph  $G$ , perfect matching  $M$ , weights  $w, k, \sigma_p, \sigma_c \in \mathbb{N}$

Question: Can  $G$  be covered by  
 $\leq \sigma_p$  alternating paths  $\nmid$   
 $\leq \sigma_c$  alternating cycles  
of total weight  $\geq k$ ?

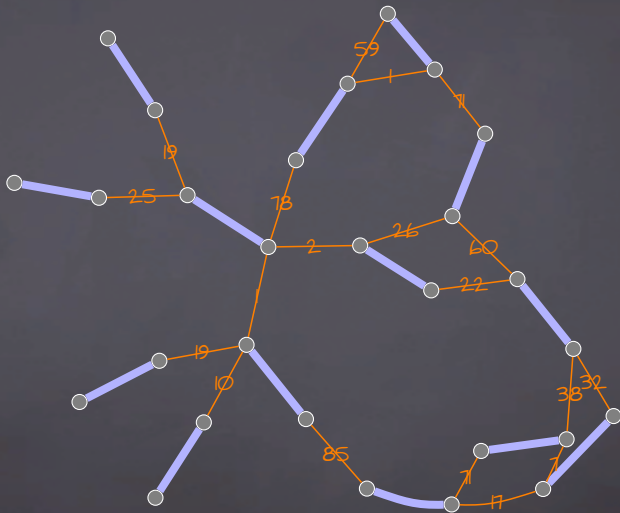
# Graph-Based Scaffolding

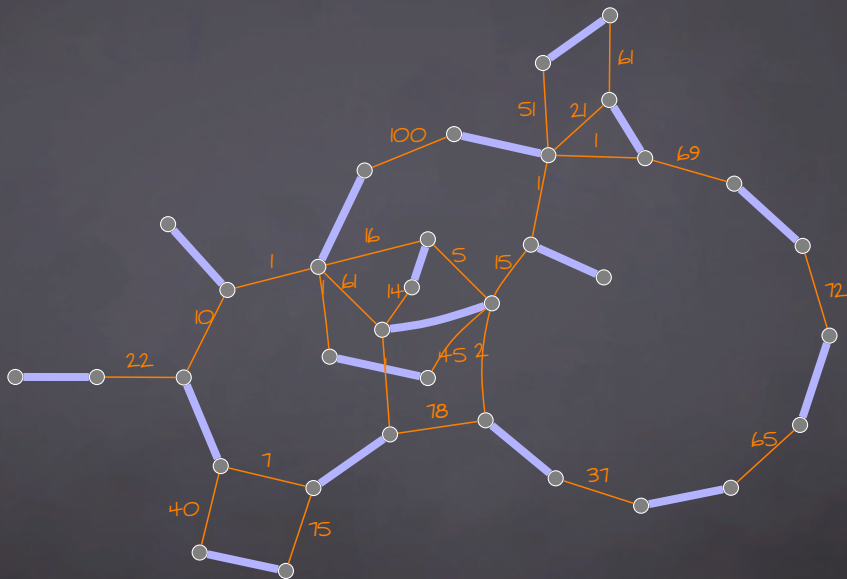


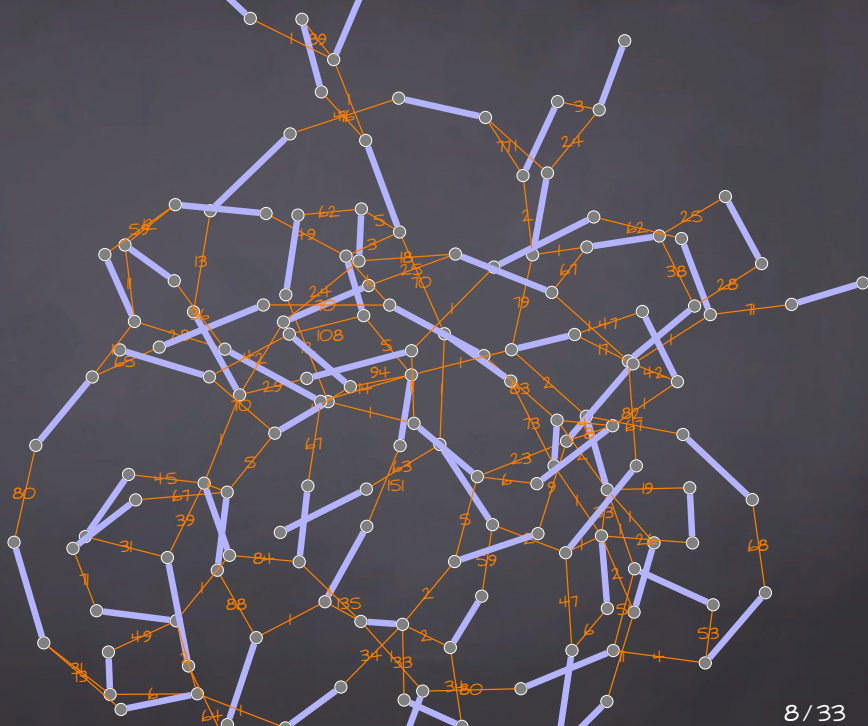
## Exact Scaffolding

Input: Graph  $G$ , perfect matching  $M$ , weights  $w, k, \sigma_p, \sigma_c \in \mathbb{N}$

Question: Can  $G$  be covered by  
 $\sigma_p$  alternating paths  $\nexists$   
 $\sigma_c$  alternating cycles  
of total weight  $\geq k$ ?





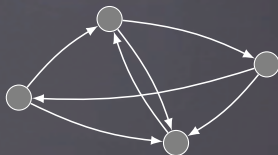
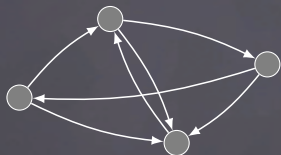


# Hardness Warm up: Hamiltonian Path

## Recall: Scaffolding

Input: Graph  $G$ , perfect matching  $M$ , weights  $w, k, \sigma_p, \sigma_c \in \mathbb{N}$

Question: Can  $G$  be covered by  $\leq \sigma_p$  alternating paths  $\neq$   
 $\leq \sigma_c$  alternating cycles of total weight  $\geq k$ ?



## Construction

Given a directed graph  $D$ .

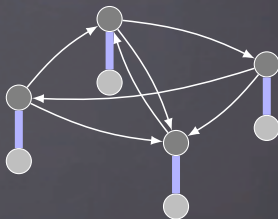
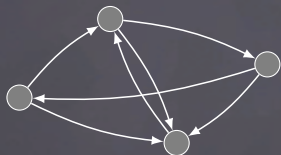
1. make a copy of  $D$

# Hardness Warm up: Hamiltonian Path

## Recall: Scaffolding

Input: Graph  $G$ , perfect matching  $M$ , weights  $w, k, \sigma_p, \sigma_c \in \mathbb{N}$

Question: Can  $G$  be covered by  $\leq \sigma_p$  alternating paths  $\neq$   
 $\leq \sigma_c$  alternating cycles of total weight  $\geq k$ ?



## Construction

Given a directed graph  $D$ .

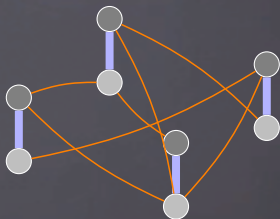
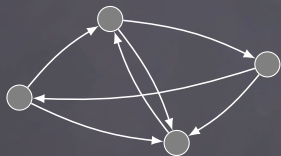
1. make a copy of  $D$
2. duplicate all vertices  $\rightsquigarrow M$

# Hardness Warm up: Hamiltonian Path

## Recall: Scaffolding

Input: Graph  $G$ , perfect matching  $M$ , weights  $w, k, \sigma_p, \sigma_c \in \mathbb{N}$

Question: Can  $G$  be covered by  $\leq \sigma_p$  alternating paths  $\nleftrightarrow$   
 $\leq \sigma_c$  alternating cycles of total weight  $\geq k$ ?



## Construction

Given a directed graph  $D$ .

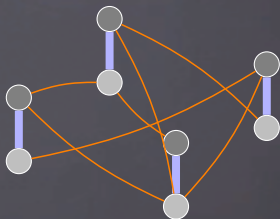
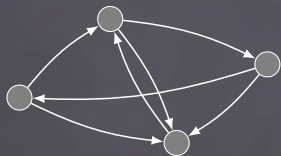
1. make a copy of  $D$
2. duplicate all vertices  $\rightsquigarrow M$
3. "slide" down all arrow tips  $\nleftrightarrow$  ignore directions

# Hardness Warm up: Hamiltonian Path

## Recall: Scaffolding

Input: Graph  $G$ , perfect matching  $M$ , weights  $w, k, \sigma_p, \sigma_c \in \mathbb{N}$

Question: Can  $G$  be covered by  $\leq \sigma_p$  alternating paths  $\neq$   
 $\leq \sigma_c$  alternating cycles of total weight  $\geq k$ ?



## Lemma

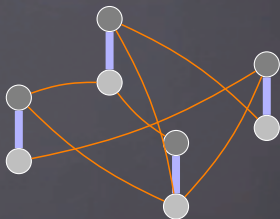
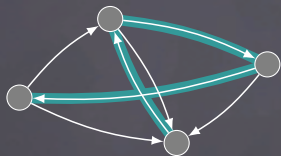
$D$  admits a directed Hamiltonian path  $\Leftrightarrow M$  can be covered with a single alternating path in  $G$ .

# Hardness Warm up: Hamiltonian Path

## Recall: Scaffolding

Input: Graph  $G$ , perfect matching  $M$ , weights  $w, k, \sigma_p, \sigma_c \in \mathbb{N}$

Question: Can  $G$  be covered by  $\leq \sigma_p$  alternating paths  $\neq$   
 $\leq \sigma_c$  alternating cycles of total weight  $\geq k$ ?



## Lemma

$D$  admits a directed Hamiltonian path  $\Leftrightarrow M$  can be covered with a single alternating path in  $G$ .

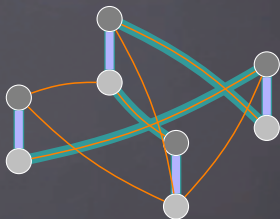
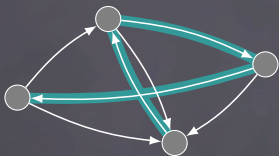
" $\Rightarrow$ ": replace each  $v$  in the Hamiltonian path by  $v_{\text{up}} \rightarrow v_{\text{low}}$ .

# Hardness Warm up: Hamiltonian Path

## Recall: Scaffolding

Input: Graph  $G$ , perfect matching  $M$ , weights  $w, k, \sigma_p, \sigma_c \in \mathbb{N}$

Question: Can  $G$  be covered by  $\leq \sigma_p$  alternating paths  $\neq$   
 $\leq \sigma_c$  alternating cycles of total weight  $\geq k$ ?



## Lemma

$D$  admits a directed Hamiltonian path  $\Leftrightarrow M$  can be covered with a single alternating path in  $G$ .

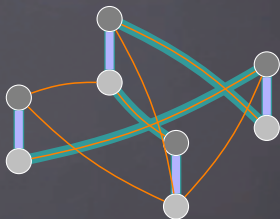
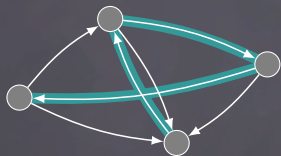
" $\Rightarrow$ ": replace each  $v$  in the Hamiltonian path by  $v_{\text{up}} \rightarrow v_{\text{low}}$ .  
alternating  $\checkmark$  covers  $M$   $\checkmark$

# Hardness Warm up: Hamiltonian Path

## Recall: Scaffolding

Input: Graph  $G$ , perfect matching  $M$ , weights  $w, k, \sigma_p, \sigma_c \in \mathbb{N}$

Question: Can  $G$  be covered by  $\leq \sigma_p$  alternating paths  $\neq$   
 $\leq \sigma_c$  alternating cycles of total weight  $\geq k$ ?



## Lemma

$D$  admits a directed Hamiltonian path  $\Leftrightarrow M$  can be covered with a single alternating path in  $G$ .

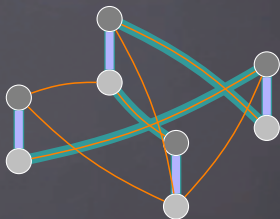
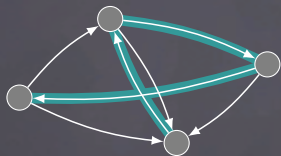
" $\Leftarrow$ ": contract each matching edge in the covering alternating path.

# Hardness Warm up: Hamiltonian Path

## Recall: Scaffolding

Input: Graph  $G$ , perfect matching  $M$ , weights  $w, k, \sigma_p, \sigma_c \in \mathbb{N}$

Question: Can  $G$  be covered by  $\leq \sigma_p$  alternating paths  $\neq$   
 $\leq \sigma_c$  alternating cycles of total weight  $\geq k$ ?



## Lemma

$D$  admits a directed Hamiltonian path  $\Leftrightarrow M$  can be covered with a single alternating path in  $G$ .

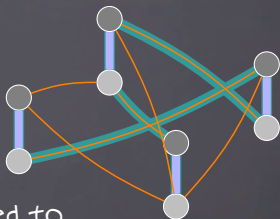
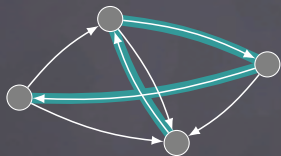
" $\Leftarrow$ ": contract each matching edge in the covering alternating path.  
hits all vertices exactly once  $\checkmark$       is valid directed path  $\checkmark$

# Hardness Warm up: Hamiltonian Path

## Recall: Scaffolding

Input: Graph  $G$ , perfect matching  $M$ , weights  $w, k, \sigma_p, \sigma_c \in \mathbb{N}$

Question: Can  $G$  be covered by  $\leq \sigma_p$  alternating paths  $\neq$   
 $\leq \sigma_c$  alternating cycles of total weight  $\geq k$ ?



## Theorem

Scaffolding is NP-hard, even restricted to

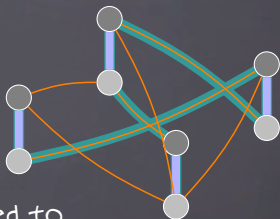
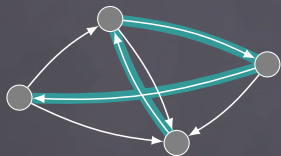
- Bipartite graphs
- $(\sigma_p, \sigma_c) \in \{(0, 1), (1, 0)\}$  and
- $w : E \rightarrow \{0\}$ .

# Hardness Warm up: Hamiltonian Path

## Recall: Scaffolding

Input: Graph  $G$ , perfect matching  $M$ , weights  $w, k, \sigma_p, \sigma_c \in \mathbb{N}$

Question: Can  $G$  be covered by  $\leq \sigma_p$  alternating paths  $\neq$   
 $\leq \sigma_c$  alternating cycles of total weight  $\geq k$ ?



## Theorem

Scaffolding is NP-hard, even restricted to

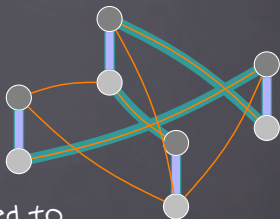
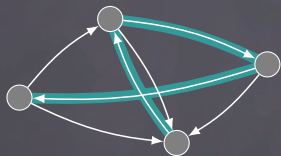
- supergraphs of bipartite graphs
- $(\sigma_p, \sigma_c) \in \{(0, 1), (1, 0)\}$  and
- $w : E \rightarrow \{0, 1\}$ .

# Hardness Warm up: Hamiltonian Path

## Recall: Scaffolding

Input: Graph  $G$ , perfect matching  $M$ , weights  $w, k, \sigma_p, \sigma_c \in \mathbb{N}$

Question: Can  $G$  be covered by  $\leq \sigma_p$  alternating paths  $\neq$   
 $\leq \sigma_c$  alternating cycles of total weight  $\geq k$ ?



## Theorem

Scaffolding is NP-hard, even restricted to

- supergraphs of bipartite graphs
- $(\sigma_p, \sigma_c) \in \{(0, 1), (1, 0)\}$  and
- $w : E \rightarrow \{0, 1\}$ .

## Corollary

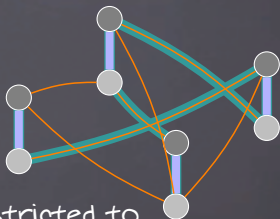
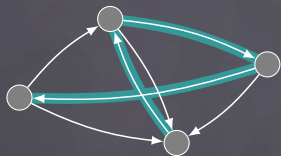
Scaffolding with 2 weights is NP-hard in any sufficiently dense graph class.

# Hardness Warm up: Hamiltonian Path

## Recall: Scaffolding

Input: Graph  $G$ , perfect matching  $M$ , weights  $w, k, \sigma_p, \sigma_c \in \mathbb{N}$

Question: Can  $G$  be covered by  $\leq \sigma_p$  alternating paths  $\neq$   
 $\leq \sigma_c$  alternating cycles of total weight  $\geq k$ ?



## Theorem

Exact Scaffolding is NP-hard, even restricted to

- supergraphs of Bipartite graphs
- $(\sigma_p, \sigma_c) \in \{(0, 1), (1, 0)\}$  and
- $w : E \rightarrow \{0, 1\}$ .

## Corollary

Exact Scaffolding with 2 weights is NP-hard in any sufficiently dense graph class.

# Scaffolding in Co-Bipartites

# Scaffolding in Co-Bipartites

Wait, what?



# Scaffolding in Co-Bipartites

Wait, what?

Recap: Corollary

Scaffolding with 2 weights is NP-hard in any sufficiently dense graph class.



# Scaffolding in Co-Bipartites

Wait, what?

Recap: Corollary

Scaffolding with 2 weights is NP-hard in any sufficiently dense graph class.

~> Unweighted!



# Scaffolding in Unweighted Co-Bipartites



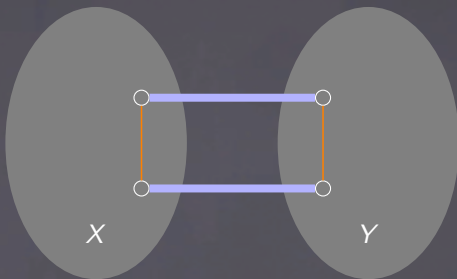
## Observation

no edges between  $X \nsubseteq Y \rightsquigarrow$  need 2 objects (paths/cycles)  
otherwise  $\rightsquigarrow$  can always cover  $G$  with 1 path

## TODO

decide if we can cover with 1 cycle

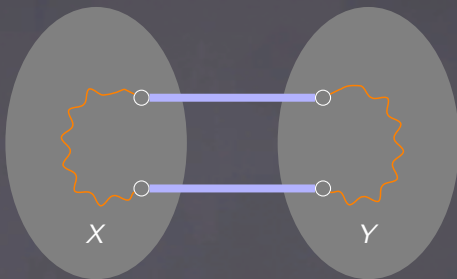
# Scaffolding in Unweighted Co-Bipartites



## Observation

- $\exists$  alternating cycle with non-matching edge  $X$
- $\rightsquigarrow$  extend to cover all  $M$  in  $G[X]$

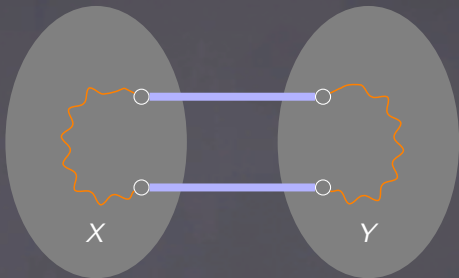
# Scaffolding in Unweighted Co-Bipartites



## Observation

$\exists$  alternating cycle with non-matching edge  $X$   
 $\rightsquigarrow$  extend to cover all  $\mathcal{M}$  in  $G[X]$

# Scaffolding in Unweighted Co-Bipartites



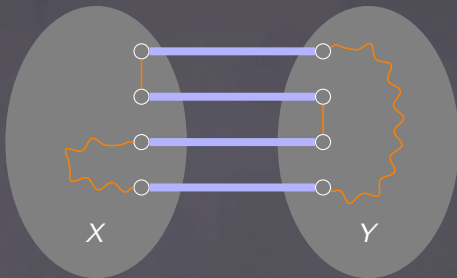
## Observation

$\exists$  alternating cycle with non-matching edge  $X$   
 $\rightsquigarrow$  extend to cover all  $M$  in  $G[X]$

## Observation

#matching edges between  $X \nsubseteq Y$  even (and  $> 0$ )  $\rightsquigarrow$  ✓

# Scaffolding in Unweighted Co-Bipartites



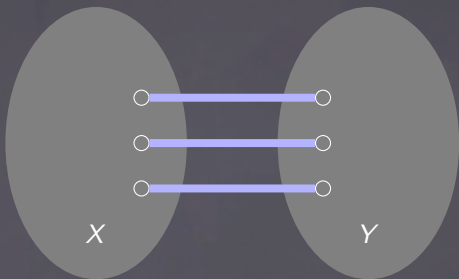
## Observation

$\exists$  alternating cycle with non-matching edge  $X$   
 $\rightsquigarrow$  extend to cover all  $M$  in  $G[X]$

## Observation

#matching edges between  $X \nsubseteq Y$  even (and  $> 0$ )  $\rightsquigarrow$  ✓

# Scaffolding in Unweighted Co-Bipartites



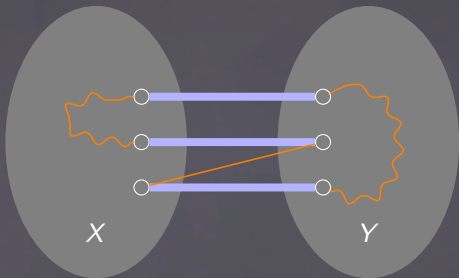
## Observation

$\exists$  alternating cycle with non-matching edge  $X$   
 $\rightsquigarrow$  extend to cover all  $M$  in  $G[X]$

## Observation

#matching edges between  $X \nleftrightarrow Y$  even (and  $> 0$ )  $\rightsquigarrow$  ✓  
#matching edges between  $X \nleftrightarrow Y$  odd

# Scaffolding in Unweighted Co-Bipartites



## Observation

$\exists$  alternating cycle with non-matching edge  $X$   
 $\rightsquigarrow$  extend to cover all  $M$  in  $G[X]$

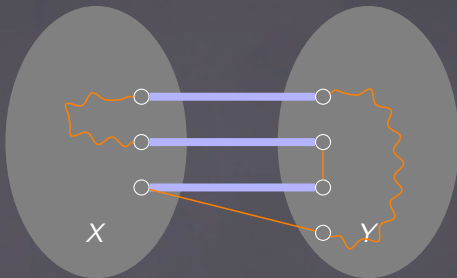
## Observation

#matching edges between  $X \nsubseteq Y$  even (and  $> 0$ )  $\rightsquigarrow$  ✓

#matching edges between  $X \nsubseteq Y$  odd

$\rightsquigarrow$  find any non-matching edge between  $X \nsubseteq Y$

# Scaffolding in Unweighted Co-Bipartites



## Observation

$\exists$  alternating cycle with non-matching edge  $X$   
 $\rightsquigarrow$  extend to cover all  $M$  in  $G[X]$

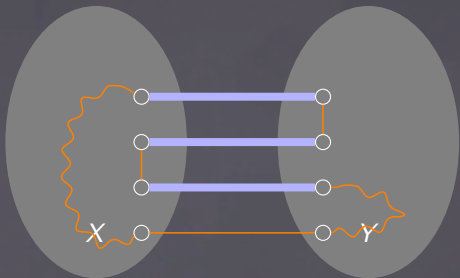
## Observation

#matching edges between  $X \nsubseteq Y$  even (and  $> 0$ )  $\rightsquigarrow$  ✓

#matching edges between  $X \nsubseteq Y$  odd

$\rightsquigarrow$  find any non-matching edge between  $X \nsubseteq Y$

# Scaffolding in Unweighted Co-Bipartites



## Observation

$\exists$  alternating cycle with non-matching edge  $X$   
 $\rightsquigarrow$  extend to cover all  $M$  in  $G[X]$

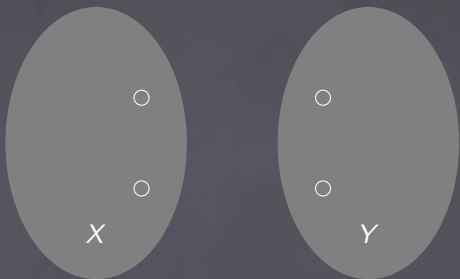
## Observation

#matching edges between  $X \nsubseteq Y$  even (and  $> 0$ )  $\rightsquigarrow$  ✓

#matching edges between  $X \nsubseteq Y$  odd

$\rightsquigarrow$  find any non-matching edge between  $X \nsubseteq Y$

# Scaffolding in Unweighted Co-Bipartites



## Observation

$\exists$  alternating cycle with non-matching edge  $X$   
 $\rightsquigarrow$  extend to cover all  $M$  in  $G[X]$

## Observation

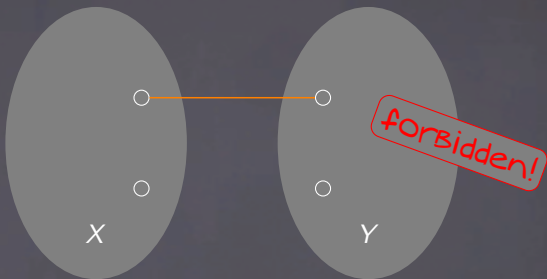
#matching edges between  $X \nleftrightarrow Y$  even (and  $> 0$ )  $\rightsquigarrow$  ✓

#matching edges between  $X \nleftrightarrow Y$  odd

$\rightsquigarrow$  find any non-matching edge between  $X \nleftrightarrow Y$

#matching edges between  $X \nleftrightarrow Y$  is 0

# Scaffolding in Unweighted Co-Bipartites



## Observation

$\exists$  alternating cycle with non-matching edge  $X$   
 $\rightsquigarrow$  extend to cover all  $M$  in  $G[X]$

## Observation

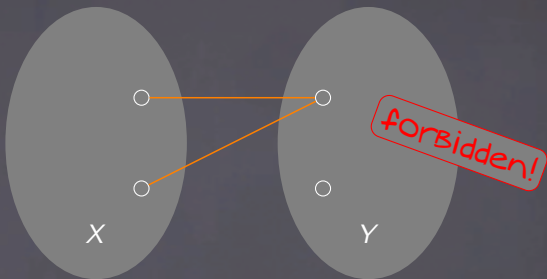
#matching edges between  $X \nleftrightarrow Y$  even (and  $> 0$ )  $\rightsquigarrow$  ✓

#matching edges between  $X \nleftrightarrow Y$  odd

$\rightsquigarrow$  find any non-matching edge between  $X \nleftrightarrow Y$

#matching edges between  $X \nleftrightarrow Y$  is 0

# Scaffolding in Unweighted Co-Bipartites



## Observation

$\exists$  alternating cycle with non-matching edge  $X$   
 $\rightsquigarrow$  extend to cover all  $M$  in  $G[X]$

## Observation

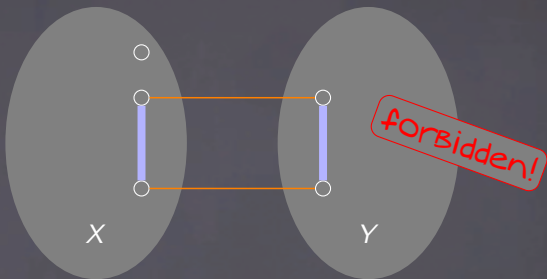
#matching edges between  $X \nsubseteq Y$  even (and  $> 0$ )  $\rightsquigarrow$  ✓

#matching edges between  $X \nsubseteq Y$  odd

$\rightsquigarrow$  find any non-matching edge between  $X \nsubseteq Y$

#matching edges between  $X \nsubseteq Y$  is 0

# Scaffolding in Unweighted Co-Bipartites



## Observation

$\exists$  alternating cycle with non-matching edge  $X$   
 $\rightsquigarrow$  extend to cover all  $M$  in  $G[X]$

## Observation

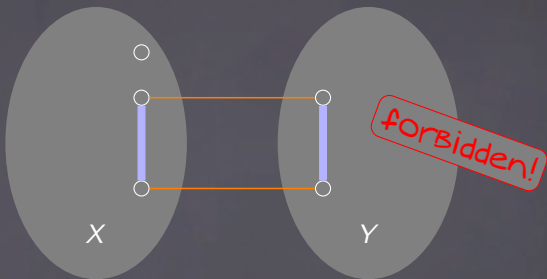
#matching edges between  $X \nsubseteq Y$  even (and  $> 0$ )  $\rightsquigarrow$  ✓

#matching edges between  $X \nsubseteq Y$  odd

$\rightsquigarrow$  find any non-matching edge between  $X \nsubseteq Y$

#matching edges between  $X \nsubseteq Y$  is 0

# Scaffolding in Unweighted Co-Bipartites



## Observation

$\exists$  alternating cycle with non-matching edge  $X$   
 $\rightsquigarrow$  extend to cover all  $M$  in  $G[X]$

## Observation

#matching edges between  $X \nsubseteq Y$  even (and  $> 0$ )  $\rightsquigarrow$   $\checkmark$

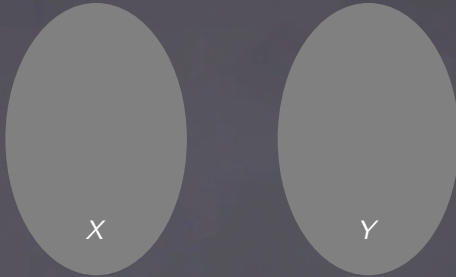
#matching edges between  $X \nsubseteq Y$  odd

$\rightsquigarrow$  find any non-matching edge between  $X \nsubseteq Y$

#matching edges between  $X \nsubseteq Y$  is 0

all other cases are  $\checkmark$  (tedious case analysis)

# Scaffolding in Unweighted Co-Bipartites



## Theorem

Scaffolding can be solved in  $O(n + m)$  time on co-Bipartite graphs

# Scaffolding in Unweighted Trees

## Observation

no alternating cycles in a tree

# Scaffolding in Unweighted Trees

## Observation

no alternating cycles in a tree

## Observation

consider a **lowest** leaf  $l$



# Scaffolding in Unweighted Trees

## Observation

no alternating cycles in a tree

## Observation

consider a **lowest** leaf  $l$

$M$  is perfect  $\rightsquigarrow l$  **matched**



# Scaffolding in Unweighted Trees

## Observation

no alternating cycles in a tree

## Observation

consider a **lowest** leaf  $l$

$M$  is perfect  $\rightsquigarrow l$  **matched**



# Scaffolding in Unweighted Trees

## Observation

no alternating cycles in a tree

## Observation

consider a **lowest** leaf  $l$

$M$  is perfect  $\rightsquigarrow l$  **matched**

parent  $p$  of  $l$  has only 1 child



# Scaffolding in Unweighted Trees

## Observation

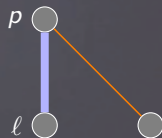
no alternating cycles in a tree

## Observation

consider a **lowest** leaf  $l$

$M$  is perfect  $\rightsquigarrow l$  **matched**

parent  $p$  of  $l$  has only 1 child



# Scaffolding in Unweighted Trees

## Observation

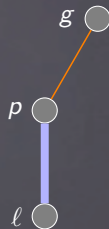
no alternating cycles in a tree

## Observation

consider a **lowest** leaf  $l$

$M$  is perfect  $\rightsquigarrow l$  **matched**

parent  $p$  of  $l$  has only 1 child



# Scaffolding in Unweighted Trees

## Observation

no alternating cycles in a tree

## Observation

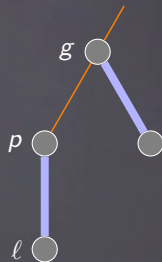
consider a **lowest** leaf  $l$

$M$  is perfect  $\rightsquigarrow l$  **matched**

parent  $p$  of  $l$  has only 1 child

## **Case 1**

parent  $g$  of  $p$  is **matched** "Below"



# Scaffolding in Unweighted Trees

## Observation

no alternating cycles in a tree

## Observation

consider a **lowest** leaf  $l$

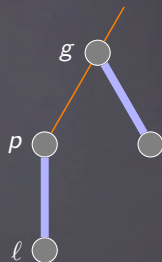
$M$  is perfect  $\rightsquigarrow l$  **matched**

parent  $p$  of  $l$  has only 1 child

## **Case 1**

parent  $g$  of  $p$  is **matched** "Below"

$\rightsquigarrow g$  is matched to a leaf  $l'$



# Scaffolding in Unweighted Trees

## Observation

no alternating cycles in a tree

## Observation

consider a **lowest** leaf  $l$

$M$  is perfect  $\leadsto l$  **matched**

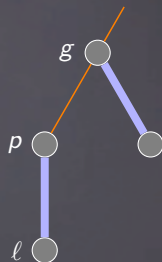
parent  $p$  of  $l$  has only 1 child

## **Case 1**

parent  $g$  of  $p$  is **matched** "Below"

$\leadsto g$  is matched to a leaf  $l'$

$\leadsto$  always take  $l-p-g-l'$



# Scaffolding in Unweighted Trees

## Observation

no alternating cycles in a tree

## Observation

consider a **lowest** leaf  $l$

$M$  is perfect  $\rightsquigarrow l$  **matched**

parent  $p$  of  $l$  has only 1 child

## Case 1

parent  $g$  of  $p$  is **matched** "Below"

$\rightsquigarrow g$  is matched to a leaf  $l'$

$\rightsquigarrow$  always take  $l-p-g-l'$

## Case 2

parent  $g$  of  $p$  is **matched** "above"



# Scaffolding in Unweighted Trees

## Observation

no alternating cycles in a tree

## Observation

consider a **lowest** leaf  $l$

$M$  is perfect  $\leadsto l$  **matched**

parent  $p$  of  $l$  has only 1 child

## Case 1

parent  $g$  of  $p$  is **matched** "Below"

$\leadsto g$  is matched to a leaf  $l'$

$\leadsto$  always take  $l-p-g-l'$

## Case 2

parent  $g$  of  $p$  is **matched** "above"

**either**  $p$  is the only child of  $g$



# Scaffolding in Unweighted Trees

## Observation

no alternating cycles in a tree

## Observation

consider a **lowest** leaf  $\ell$

$M$  is perfect  $\leadsto \ell$  **matched**

parent  $p$  of  $\ell$  has only 1 child



## Case 1

parent  $g$  of  $p$  is **matched** "Below"

$\leadsto g$  is matched to a leaf  $\ell'$

$\leadsto$  always take  $\ell - p - g - \ell'$

## Case 2

parent  $g$  of  $p$  is **matched** "above"

**either**  $p$  is the only child of  $g \leadsto$  delete  $\ell \neq g$  and reduce  $k$

# Scaffolding in Unweighted Trees

## Observation

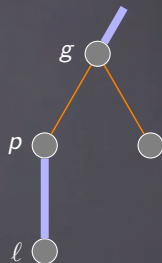
no alternating cycles in a tree

## Observation

consider a **lowest** leaf  $l$

$M$  is perfect  $\leadsto l$  **matched**

parent  $p$  of  $l$  has only 1 child



## Case 1

parent  $g$  of  $p$  is **matched** "Below"

$\leadsto g$  is matched to a leaf  $l'$

$\leadsto$  always take  $l-p-g-l'$

## Case 2

parent  $g$  of  $p$  is **matched** "above"

**either**  $p$  is the only child of  $g \leadsto$  delete  $l \notin g$  and reduce  $k$

**or**  $g$  has another child  $u$

# Scaffolding in Unweighted Trees

## Observation

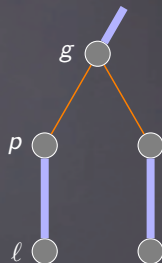
no alternating cycles in a tree

## Observation

consider a **lowest** leaf  $l$

$M$  is perfect  $\leadsto l$  **matched**

parent  $p$  of  $l$  has only 1 child



## Case 1

parent  $g$  of  $p$  is **matched** "Below"

$\leadsto g$  is matched to a leaf  $l'$

$\leadsto$  always take  $l-p-g-l'$

## Case 2

parent  $g$  of  $p$  is **matched** "above"

**either**  $p$  is the only child of  $g \leadsto$  delete  $l \neq g$  and reduce  $k$

**or**  $g$  has another child  $u \leadsto u$  matched "Below"

## Scaffolding in Unweighted Trees

## Observation

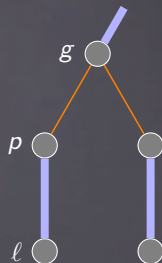
no alternating cycles in a tree

## Observation

consider a lowest leaf  $l$

$M$  is perfect  $\leadsto l$  matched

parent  $p$  of  $l$  has only 1 child



## Case I

parent  $g$  of  $p$  is matched "Below"

$\leadsto g$  is matched to a leaf  $l'$

$\leadsto$  always take  $\ell - p - g - \ell'$

## Case 2

parent  $g$  of  $p$  is matched "above"

either  $p$  is the only child of  $g \rightsquigarrow$  delete  $l \notin g$  and reduce  $k$

or  $g$  has another child  $u \rightsquigarrow u$  matched "Below"  $\rightsquigarrow \exists$  "clone" of  $g-p-l$

# Scaffolding in Unweighted Trees

## Observation

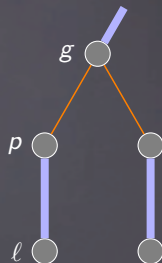
no alternating cycles in a tree

## Observation

consider a **lowest** leaf  $l$

$M$  is perfect  $\leadsto l$  **matched**

parent  $p$  of  $l$  has only 1 child



## Case 1

parent  $g$  of  $p$  is **matched** "Below"

$\leadsto g$  is matched to a leaf  $l'$

$\leadsto$  always take  $l-p-g-l'$

## Case 2

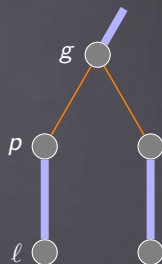
parent  $g$  of  $p$  is **matched** "above"

**either**  $p$  is the only child of  $g \leadsto$  delete  $l \neq g$  and reduce  $k$

**or**  $g$  has another child  $u \leadsto u$  matched "Below"  $\leadsto \exists$  "clone" of  $g-p-l$

$\leadsto$  take  $p-l$

# Scaffolding in Unweighted Trees



## Theorem

Scaffolding can be solved in  $O(n)$  time on unweighted trees

# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used below  $v$
- $v$  incident with non-matching?

# Scaffolding in Weighted Trees

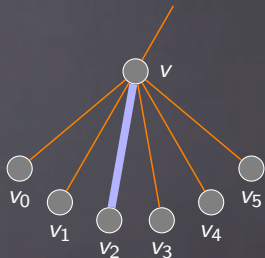
## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used Below  $v$
- $v$  incident with non-matching?

## Semantics

$[p, x]_v = \text{max weight collected below } v \text{ with } p \text{ finished paths "under } x\text{"}$



# Scaffolding in Weighted Trees

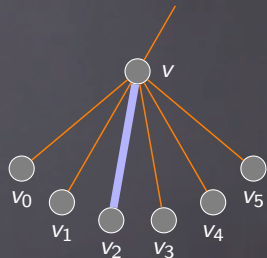
## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used below  $v$
- $v$  incident with non-matching?

## Semantics

$[p, x]_v = \max$  weight collected below  $v$  with  $p$  finished paths "under  $x$ "



## Recurrence

Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .

$$[p, x]_v := \max_{\substack{p_1, p_2, \dots, p_c \\ \sum p_i = p}} \sum_{1 \leq i \leq c} \max_{x \in \{\checkmark, \times\}} [p_i, x]_{v_i}$$

# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used below  $v$
- $v$  incident with non-matching?

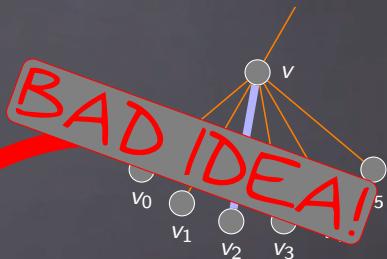
## Semantics

$[p, x]_v = \max$  weight collected below  $v$  with  $p$  finished paths "under  $x$ "

## Recurrence

Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .

$$[p, x]_v := \max_{\substack{p_1, p_2, \dots, p_c \\ \sum p_i = p}} \sum_{1 \leq i \leq c} \max_{x \in \{\checkmark, \times\}} [p_i, x]_{v_i}$$



# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used below  $v$
- $v$  incident with non-matching?

## Semantics

$[j, p, x]_v = \text{max weight collected}$

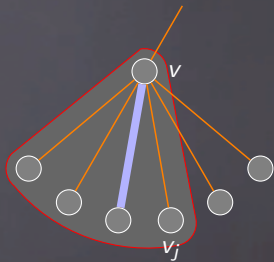
Below  $v$  with  $p$  finished paths

"under  $x$ " up to  $v_j$

(abbrev: last child  $\rightsquigarrow [p, x]_v$ )

## Recurrence

Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .



# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used Below  $v$
- $v$  incident with non-matching?

## Semantics

$[j, p, x]_v = \text{max weight collected}$

Below  $v$  with  $p$  finished paths

"under  $x$ " up to  $v_j$

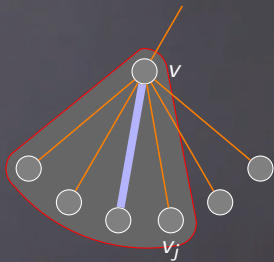
(abbrev: last child  $\rightsquigarrow [p, x]_v$ )

## Recurrence

Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .

$$[0, 0, \text{X}]_v := 0$$

$$[j, p, x]_v := \max_{p_j \leq p} \left\{ \begin{array}{l} \end{array} \right.$$



# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used Below  $v$
- $v$  incident with non-matching?

## Semantics

$[j, p, x]_v = \text{max weight collected}$

Below  $v$  with  $p$  finished paths

"under  $x$ " up to  $v_j$

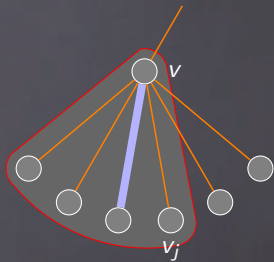
(abbrev: last child  $\rightsquigarrow [p, x]_v$ )

## Recurrence

Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .

$$[0, 0, \text{X}]_v := 0$$

$$[j, p, x]_v := \max_{p_j \leq p} \begin{cases} \max\{[p_j, \text{X}]_{v_j}, [p_j, \text{X}]_{v_j}\} + [j-1, p-p_j, x]_v & \text{if } vv_j \notin \mathcal{M} \\ \dots \end{cases}$$



# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used below  $v$
- $v$  incident with non-matching?

## Semantics

$[j, p, x]_v = \text{max weight collected}$

Below  $v$  with  $p$  finished paths

"under  $x$ " up to  $v_j$

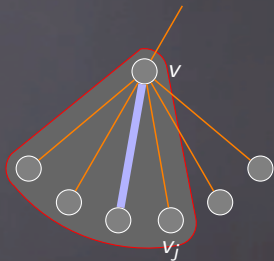
(abbrev: last child  $\rightsquigarrow [p, x]_v$ )

## Recurrence

Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .

$$[0, 0, \text{✗}]_v := 0$$

$$[j, p, x]_v := \max_{p_j \leq p} \begin{cases} \max\{[p_j, \text{✓}]_{v_j}, [p_j, \text{✗}]_{v_j}\} + [j-1, p-p_j, x]_v & \text{if } vv_j \notin \mathcal{M} \\ \omega(vv_j) + [p_j+1, \text{✗}]_{v_j} + [j-1, p-p_j, \text{✗}]_v & \text{if } x = \text{✓} \nmid vv_j \notin \mathcal{M} \end{cases}$$



# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used below  $v$
- $v$  incident with non-matching?

## Semantics

$[j, p, x]_v = \text{max weight collected}$

Below  $v$  with  $p$  finished paths

"under  $x$ " up to  $v_j$

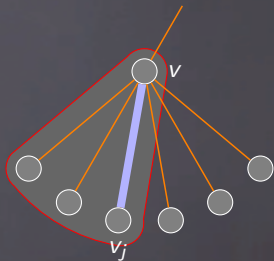
(abbrev: last child  $\rightsquigarrow [p, x]_v$ )

## Recurrence

Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .

$$[0, 0, \text{✗}]_v := 0$$

$$[j, p, x]_v := \max_{p_j \leq p} \begin{cases} \max\{[p_j, \text{✓}]_{v_j}, [p_j, \text{✗}]_{v_j}\} + [j-1, p-p_j, x]_v & \text{if } vv_j \notin \mathcal{M} \\ \omega(vv_j) + [p_j+1, \text{✗}]_{v_j} + [j-1, p-p_j, \text{✗}]_v & \text{if } x = \text{✓} \nmid vv_j \notin \mathcal{M} \\ \left\{ \begin{array}{l} [p_j-1, \text{✗}]_{v_j} \\ [p_j-1, \text{✓}]_{v_j} \end{array} \right\} + [j-1, p-p_j, x]_v & \text{if } vv_j \in \mathcal{M} \end{cases}$$



# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used Below  $v$
- $v$  incident with non-matching?

## Semantics

$[j, p, x]_v = \text{max weight collected below } v \text{ with } p \text{ finished paths}$

"under  $x$ " up to  $v_j$

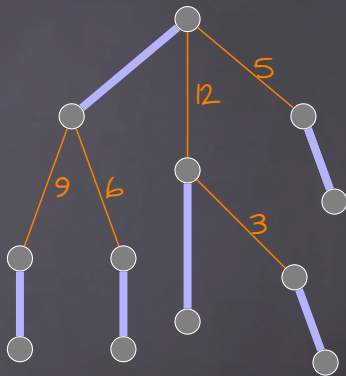
(abbrev: last child  $\rightsquigarrow [p, x]_v$ )

## Recurrence

Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .

$$[0, 0, \text{✗}]_v := 0$$

$$[j, p, x]_v := \max_{p_j \leq p} \begin{cases} \max\{[p_j, \text{✓}]_{v_j}, [p_j, \text{✗}]_{v_j}\} + [j-1, p-p_j, x]_v & \text{if } vv_j \notin \mathcal{M} \\ \omega(vv_j) + [p_j+1, \text{✗}]_{v_j} + [j-1, p-p_j, \text{✗}]_v & \text{if } x = \text{✓} \notin vv_j \notin \mathcal{M} \\ \left\{ \begin{array}{l} [p_j-1, \text{✗}]_{v_j} \\ [p_j-1, \text{✓}]_{v_j} \end{array} \right\} + [j-1, p-p_j, x]_v & \text{if } vv_j \in \mathcal{M} \end{cases}$$



# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used Below  $v$
- $v$  incident with non-matching?

## Semantics

$[j, p, x]_v = \text{max weight collected below } v \text{ with } p \text{ finished paths}$

"under  $x$ " up to  $v_j$

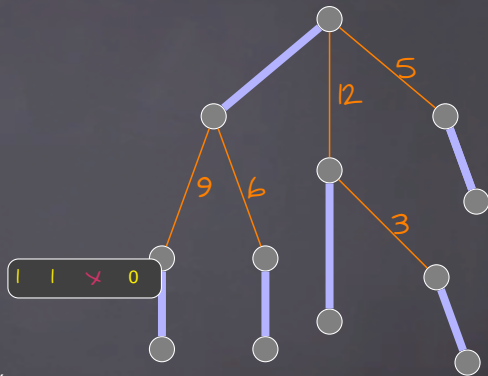
(abbrev: last child  $\rightsquigarrow [p, x]_v$ )

## Recurrence

Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .

$$[0, 0, \text{✗}]_v := 0$$

$$[j, p, x]_v := \max_{p_j \leq p} \begin{cases} \max\{[p_j, \text{✓}]_{v_j}, [p_j, \text{✗}]_{v_j}\} + [j-1, p-p_j, x]_v & \text{if } vv_j \notin \mathcal{M} \\ \omega(vv_j) + [p_j+1, \text{✗}]_{v_j} + [j-1, p-p_j, \text{✗}]_v & \text{if } x = \text{✓} \neq vv_j \notin \mathcal{M} \\ \left\{ \begin{array}{l} [p_j-1, \text{✗}]_{v_j} \\ [p_j-1, \text{✓}]_{v_j} \end{array} \right\} + [j-1, p-p_j, x]_v & \text{if } vv_j \in \mathcal{M} \end{cases}$$



# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used below  $v$
- $v$  incident with non-matching?

## Semantics

$[j, p, x]_v = \max$  weight collected below  $v$  with  $p$  finished paths

"under  $x$ " up to  $v_j$

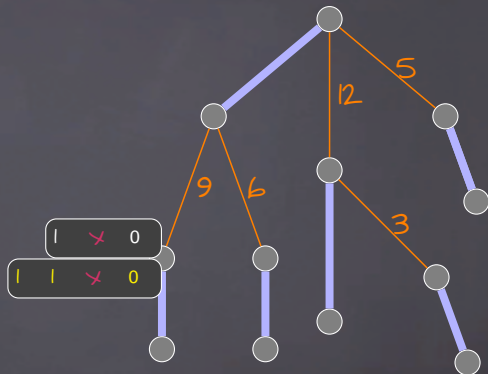
(abbrev: last child  $\rightsquigarrow [p, x]_v$ )

## Recurrence

Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .

$$[0, 0, \text{X}]_v := 0$$

$$[j, p, x]_v := \max_{p_j \leq p} \begin{cases} \max\{[p_j, \checkmark]_{v_j}, [p_j, \text{X}]_{v_j}\} + [j-1, p-p_j, x]_v & \text{if } vv_j \notin \mathcal{M} \\ \omega(vv_j) + [p_j+1, \text{X}]_{v_j} + [j-1, p-p_j, \text{X}]_v & \text{if } x = \checkmark \neq vv_j \notin \mathcal{M} \\ \left\{ \begin{array}{l} [p_j-1, \text{X}]_{v_j} \\ [p_j-1, \checkmark]_{v_j} \end{array} \right\} + [j-1, p-p_j, x]_v & \text{if } vv_j \in \mathcal{M} \end{cases}$$



# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used below  $v$
- $v$  incident with non-matching?

## Semantics

$[j, p, x]_v = \max$  weight collected below  $v$  with  $p$  finished paths

"under  $x$ " up to  $v_j$

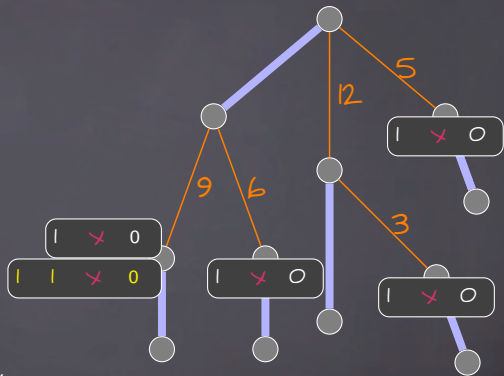
(abbrev: last child  $\rightsquigarrow [p, x]_v$ )

## Recurrence

Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .

$$[0, 0, \text{X}]_v := 0$$

$$[j, p, x]_v := \max_{p_j \leq p} \begin{cases} \max\{[p_j, \checkmark]_{v_j}, [p_j, \text{X}]_{v_j}\} + [j-1, p-p_j, x]_v & \text{if } vv_j \notin \mathcal{M} \\ \omega(vv_j) + [p_j+1, \text{X}]_{v_j} + [j-1, p-p_j, \text{X}]_v & \text{if } x = \checkmark \neq vv_j \notin \mathcal{M} \\ \left\{ \begin{array}{l} [p_j-1, \text{X}]_{v_j} \\ [p_j-1, \checkmark]_{v_j} \end{array} \right\} + [j-1, p-p_j, x]_v & \text{if } vv_j \in \mathcal{M} \end{cases}$$



# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used Below  $v$
- $v$  incident with non-matching?

## Semantics

$[j, p, x]_v = \max$  weight collected Below  $v$  with  $p$  finished paths

"under  $x$ " up to  $v_j$

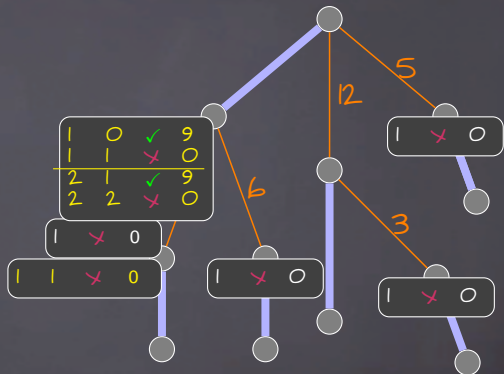
(abbrev: last child  $\rightsquigarrow [p, x]_v$ )

## Recurrence

Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .

$$[0, 0, \text{X}]_v := 0$$

$$[j, p, x]_v := \max_{p_j \leq p} \begin{cases} \max\{[p_j, \checkmark]_{v_j}, [p_j, \text{X}]_{v_j}\} + [j-1, p-p_j, x]_v & \text{if } vv_j \notin \mathcal{M} \\ \omega(vv_j) + [p_j+1, \text{X}]_{v_j} + [j-1, p-p_j, \text{X}]_v & \text{if } x = \checkmark \neq vv_j \notin \mathcal{M} \\ \left\{ \begin{array}{l} [p_j-1, \text{X}]_{v_j} \\ [p_j-1, \checkmark]_{v_j} \end{array} \right\} + [j-1, p-p_j, x]_v & \text{if } vv_j \in \mathcal{M} \end{cases}$$



# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used below  $v$
- $v$  incident with non-matching?

## Semantics

$[j, p, x]_v = \max$  weight collected below  $v$  with  $p$  finished paths

"under  $x$ " up to  $v_j$

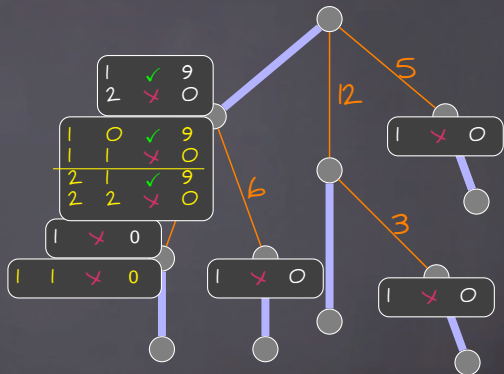
(abbrev: last child  $\rightsquigarrow [p, x]_v$ )

## Recurrence

Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .

$$[0, 0, \text{X}]_v := 0$$

$$[j, p, x]_v := \max_{p_j \leq p} \begin{cases} \max\{[p_j, \checkmark]_{v_j}, [p_j, \text{X}]_{v_j}\} + [j-1, p-p_j, x]_v & \text{if } vv_j \notin \mathcal{M} \\ \omega(vv_j) + [p_j+1, \text{X}]_{v_j} + [j-1, p-p_j, \text{X}]_v & \text{if } x = \checkmark \neq vv_j \notin \mathcal{M} \\ \left\{ \begin{array}{l} [p_j-1, \text{X}]_{v_j} \\ [p_j-1, \checkmark]_{v_j} \end{array} \right\} + [j-1, p-p_j, x]_v & \text{if } vv_j \in \mathcal{M} \end{cases}$$



# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used Below  $v$
- $v$  incident with non-matching?

## Semantics

$[j, p, x]_v = \max$  weight collected Below  $v$  with  $p$  finished paths

"under  $x$ " up to  $v_j$

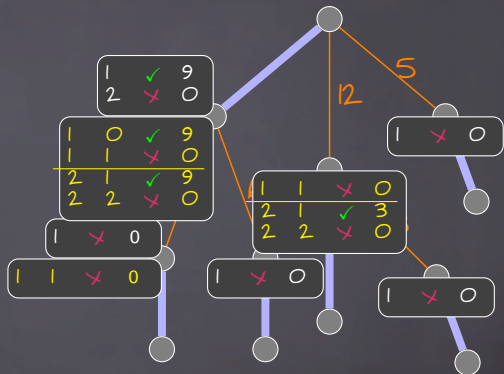
(abbrev: last child  $\rightsquigarrow [p, x]_v$ )

## Recurrence

Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .

$$[0, 0, \text{X}]_v := 0$$

$$[j, p, x]_v := \max_{p_j \leq p} \begin{cases} \max\{[p_j, \checkmark]_{v_j}, [p_j, \text{X}]_{v_j}\} + [j-1, p-p_j, x]_v & \text{if } vv_j \notin \mathcal{M} \\ \omega(vv_j) + [p_j+1, \text{X}]_{v_j} + [j-1, p-p_j, \text{X}]_v & \text{if } x = \checkmark \neq vv_j \notin \mathcal{M} \\ \left\{ \begin{array}{l} [p_j-1, \text{X}]_{v_j} \\ [p_j-1, \checkmark]_{v_j} \end{array} \right\} + [j-1, p-p_j, x]_v & \text{if } vv_j \in \mathcal{M} \end{cases}$$



# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used Below  $v$
- $v$  incident with non-matching?

## Semantics

$[j, p, x]_v = \max$  weight collected Below  $v$  with  $p$  finished paths

"under  $x$ " up to  $v_j$

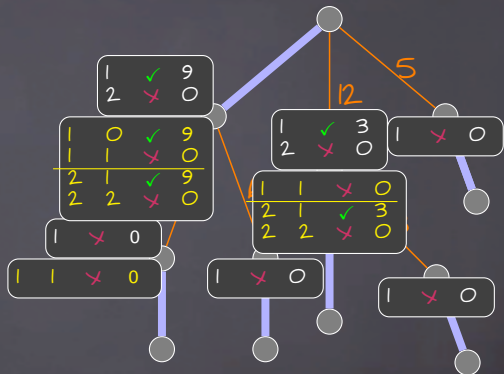
(abbrev: last child  $\rightsquigarrow [p, x]_v$ )

## Recurrence

Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .

$$[0, 0, \text{X}]_v := 0$$

$$[j, p, x]_v := \max_{p_j \leq p} \begin{cases} \max\{[p_j, \checkmark]_{v_j}, [p_j, \text{X}]_{v_j}\} + [j-1, p-p_j, x]_v & \text{if } vv_j \notin \mathcal{M} \\ \omega(vv_j) + [p_j+1, \text{X}]_{v_j} + [j-1, p-p_j, \text{X}]_v & \text{if } x = \checkmark \neq vv_j \notin \mathcal{M} \\ \left\{ \begin{array}{l} [p_j-1, \text{X}]_{v_j} \\ [p_j-1, \checkmark]_{v_j} \end{array} \right\} + [j-1, p-p_j, x]_v & \text{if } vv_j \in \mathcal{M} \end{cases}$$



# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used Below  $v$
- $v$  incident with non-matching?

## Semantics

$[j, p, x]_v = \max$  weight collected Below  $v$  with  $p$  finished paths

"under  $x$ " up to  $v_j$

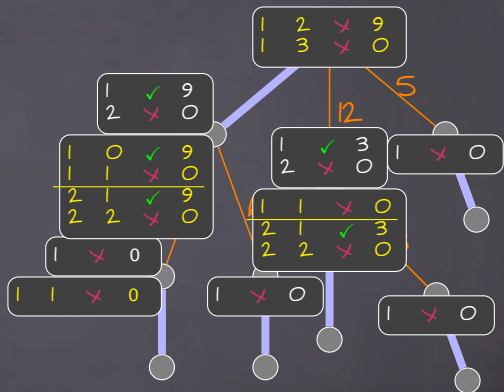
(abbrev: last child  $\rightsquigarrow [p, x]_v$ )

## Recurrence

Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .

$$[0, 0, \text{X}]_v := 0$$

$$[j, p, x]_v := \max_{p_j \leq p} \begin{cases} \max\{[p_j, \checkmark]_{v_j}, [p_j, \text{X}]_{v_j}\} + [j-1, p-p_j, x]_v & \text{if } vv_j \notin \mathcal{M} \\ \omega(vv_j) + [p_j+1, \text{X}]_{v_j} + [j-1, p-p_j, \text{X}]_v & \text{if } x = \checkmark \neq vv_j \notin \mathcal{M} \\ \left\{ \begin{array}{l} [p_j-1, \text{X}]_{v_j} \\ [p_j-1, \checkmark]_{v_j} \end{array} \right\} + [j-1, p-p_j, x]_v & \text{if } vv_j \in \mathcal{M} \end{cases}$$



# Scaffolding in Weighted Trees

## Dynamic Programming Idea

Bottom-up traversal; in each vertex  $v$ , need to remember:

- #paths used Below  $v$
- $v$  incident with non-matching?

## Semantics

$[j, p, x]_v = \max$  weight collected Below  $v$  with  $p$  finished paths

"under  $x$ " up to  $v_j$

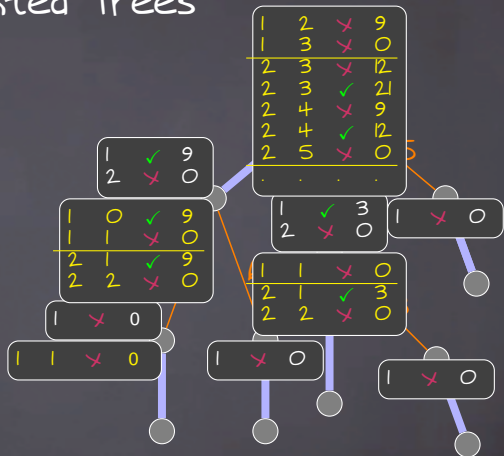
(abbrev: last child  $\rightsquigarrow [p, x]_v$ )

## Recurrence

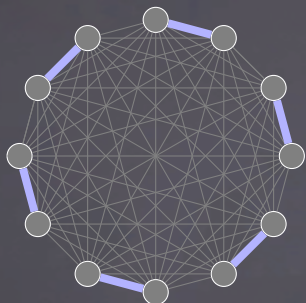
Let  $v_1, v_2, \dots, v_c$  be the children of  $v$ .

$$[0, 0, \text{X}]_v := 0$$

$$[j, p, x]_v := \max_{p_j \leq p} \begin{cases} \max\{[p_j, \checkmark]_{v_j}, [p_j, \text{X}]_{v_j}\} + [j-1, p-p_j, x]_v & \text{if } vv_j \notin \mathcal{M} \\ \omega(vv_j) + [p_j+1, \text{X}]_{v_j} + [j-1, p-p_j, \text{X}]_v & \text{if } x = \checkmark \neq vv_j \notin \mathcal{M} \\ \left\{ \begin{array}{l} [p_j-1, \text{X}]_{v_j} \\ [p_j-1, \checkmark]_{v_j} \end{array} \right\} + [j-1, p-p_j, x]_v & \text{if } vv_j \in \mathcal{M} \end{cases}$$



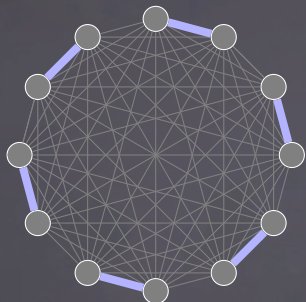
# 3-Approximation in Dense Graphs



$$\sigma_p = 1, \sigma_c = 1?$$

Approximate Scaffolding

### 3-Approximation in Dense Graphs

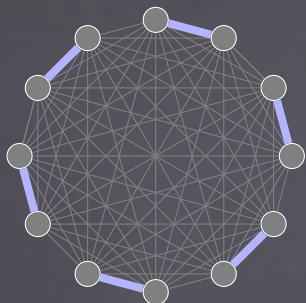


$$\sigma_p = 1, \sigma_c = 1?$$

Approximate *Scaffolding*

1. sort all edges by weight

### 3-Approximation in Dense Graphs

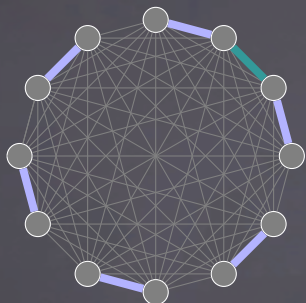


$$\sigma_p = 1, \sigma_c = 1?$$

#### Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

### 3-Approximation in Dense Graphs

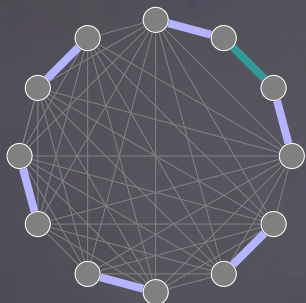


$$\sigma_p = 1, \sigma_c = 1?$$

#### Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

### 3-Approximation in Dense Graphs

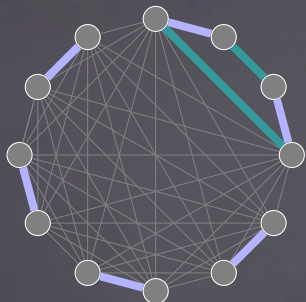


$$\sigma_p = 1, \sigma_c = 1?$$

#### Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

### 3-Approximation in Dense Graphs

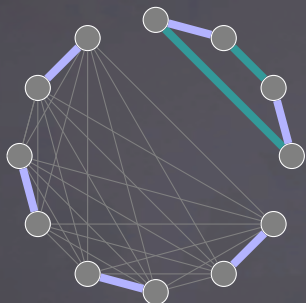


$$\sigma_p = 1, \sigma_c = 1?$$

#### Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

### 3-Approximation in Dense Graphs

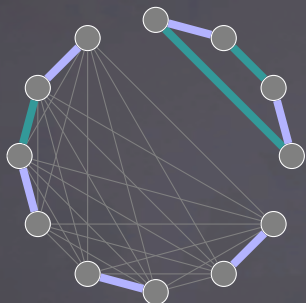


$$\sigma_p = 1, \sigma_c = 1?$$

#### Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

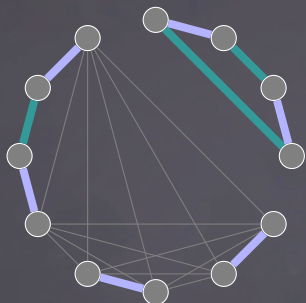
### 3-Approximation in Dense Graphs



#### Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

### 3-Approximation in Dense Graphs

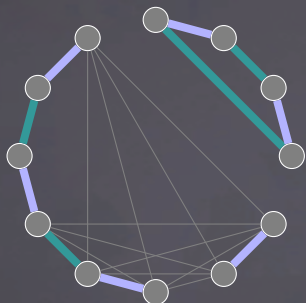


$$\sigma_p = 1, \sigma_c = 1?$$

#### Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

### 3-Approximation in Dense Graphs

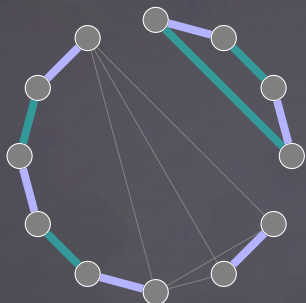


$$\sigma_p = 1, \sigma_c = 1?$$

#### Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

### 3-Approximation in Dense Graphs

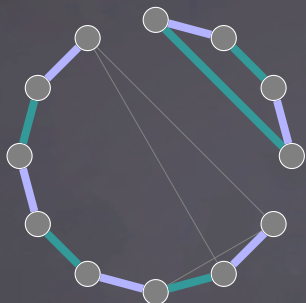


$\sigma_p = 1, \sigma_c = 1?$

#### Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

### 3-Approximation in Dense Graphs

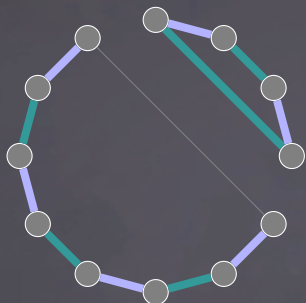


$$\sigma_p = 1, \sigma_c = 1?$$

#### Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

### 3-Approximation in Dense Graphs

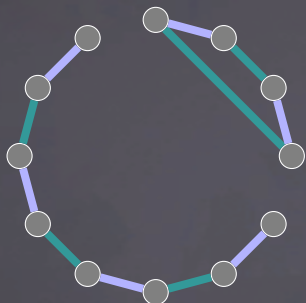


$\sigma_p = 1, \sigma_c = 1?$

#### Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

### 3-Approximation in Dense Graphs



$$\sigma_p = 1, \sigma_c = 1?$$

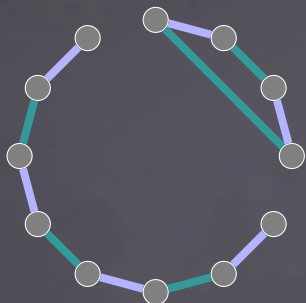
#### Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

#### Proof

Result  $S^*$  is a valid solution ✓

# 3-Approximation in Dense Graphs



$$\sigma_p = 1, \sigma_c = 1?$$

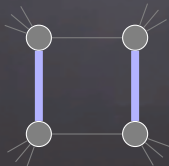
## Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

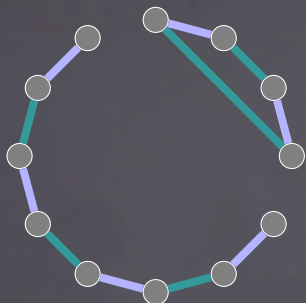
## Proof

Result  $S^*$  is a valid solution ✓

Note: taking an edge forbids  $\leq 3$  OPT edges



# 3-Approximation in Dense Graphs



$$\sigma_p = 1, \sigma_c = 1?$$

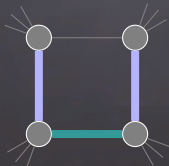
## Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

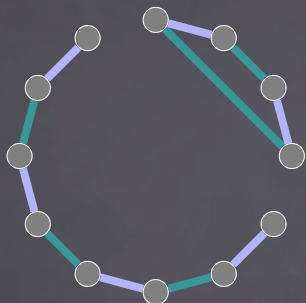
## Proof

Result  $S^*$  is a valid solution ✓

Note: taking an edge forbids  $\leq 3$  OPT edges



# 3-Approximation in Dense Graphs



$$\sigma_p = 1, \sigma_c = 1?$$

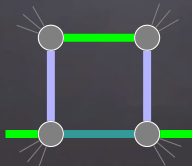
## Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

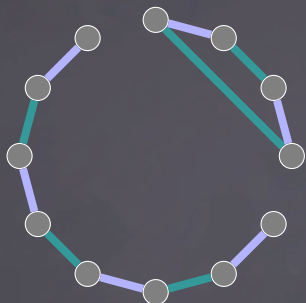
## Proof

Result  $S^*$  is a valid solution ✓

Note: taking an edge forbids  $\leq 3$  OPT edges



### 3-Approximation in Dense Graphs



$$\sigma_p = 1, \sigma_c = 1?$$

#### Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

#### Proof

Result  $S^*$  is a valid solution ✓

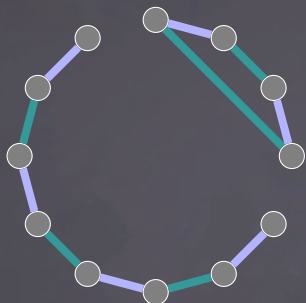
Note: taking an edge forbids  $\leq 3$  OPT edges

⇒ mark the  $\leq 3$  OPT-edges when taking an edge  $e$

⇒  $e$  is heaviest among them

⇒  $3w(S^*) \geq OPT$

### 3-Approximation in Dense Graphs



$$\sigma_p = 1, \sigma_c = 1?$$

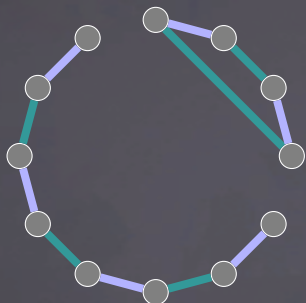
#### Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

#### Theorem

Scaffolding in complete graphs can be 3-approximated in  $O(|V| \log |V|)$  time.

### 3-Approximation in Dense Graphs



$$\sigma_p = 1, \sigma_c = 1?$$

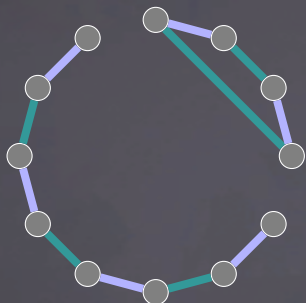
#### Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

#### Theorem

Scaffolding in complete (Bipartite) graphs can be 3-approximated in  $O(|V| \log |V|)$  time.

### 3-Approximation in Dense Graphs



$$\sigma_p = 1, \sigma_c = 1?$$

#### Approximate Scaffolding

1. sort all edges by weight
2. repeatedly take heaviest edge, if possible

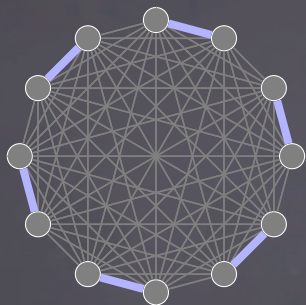
#### Theorem

Scaffolding in complete (Bipartite) graphs can be 3-approximated in  $O(|V| \log |V|)$  time.

#### Remark

For Exact Scaffolding, we have to keep an eye on the number of components too.

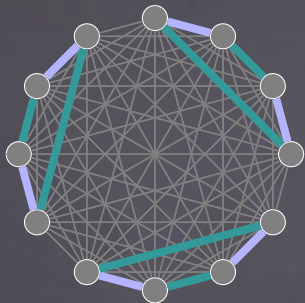
## 2-Approximation in Dense Graphs



$$\sigma_p = 1, \sigma_c = 1?$$

Approximate **Exact Scaffolding**

## 2-Approximation in Dense Graphs



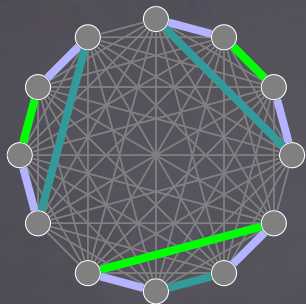
$$\sigma_p = 1, \sigma_c = 1?$$

### Approximate Exact Scaffolding

1. compute max-weight perfect matching  $S$

$\leadsto S \cup M$  is collection of cycles

## 2-Approximation in Dense Graphs

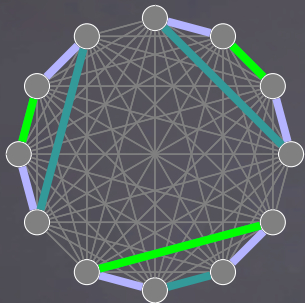


$$\sigma_p = 1, \sigma_c = 1?$$

### Approximate Exact Scaffolding

1. compute max-weight perfect matching  $S$   
 $\leadsto S \cup M$  is collection of cycles
2. "fix" all but lightest edge per cycle

## 2-Approximation in Dense Graphs

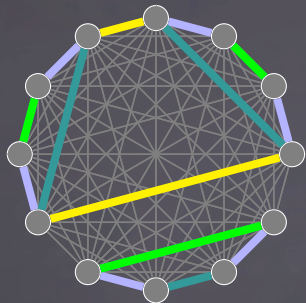


$$\sigma_p = 1, \sigma_c = 1?$$

### Approximate Exact Scaffolding

1. compute max-weight perfect matching  $S$   
 $\leadsto S \cup M$  is collection of cycles
2. "fix" all but lightest edge per cycle
3. repeatedly flip any lightest non-fix 4-cycle intersecting 2 cycles until at most  $\sigma_c + \sigma_p$  cycles remain

## 2-Approximation in Dense Graphs

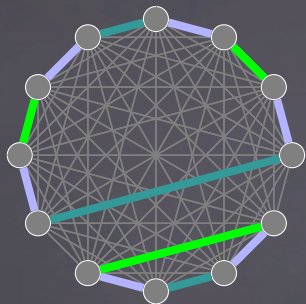


$$\sigma_p = 1, \sigma_c = 1?$$

### Approximate Exact Scaffolding

1. compute max-weight perfect matching  $S$   
 $\leadsto S \cup M$  is collection of cycles
2. "fix" all but lightest edge per cycle
3. repeatedly flip any lightest non-fix 4-cycle intersecting 2 cycles until at most  $\sigma_c + \sigma_p$  cycles remain

## 2-Approximation in Dense Graphs

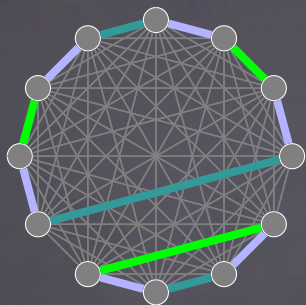


$$\sigma_p = 1, \sigma_c = 1?$$

### Approximate Exact Scaffolding

1. compute max-weight perfect matching  $S$   
 $\leadsto S \cup M$  is collection of cycles
2. "fix" all but lightest edge per cycle
3. repeatedly flip any lightest non-fix 4-cycle intersecting 2 cycles until at most  $\sigma_c + \sigma_p$  cycles remain

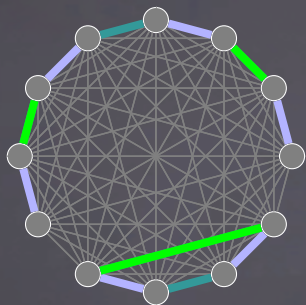
## 2-Approximation in Dense Graphs



### Approximate Exact Scaffolding

1. compute max-weight perfect matching  $S$   
 $\leadsto S \cup M$  is collection of cycles
2. "fix" all but lightest edge per cycle
3. repeatedly flip any lightest non-fix 4-cycle intersecting 2 cycles until at most  $\sigma_c + \sigma_p$  cycles remain
4. repeatedly remove lightest non-fix cycle-edge until at most  $\sigma_c$  cycles remain

## 2-Approximation in Dense Graphs

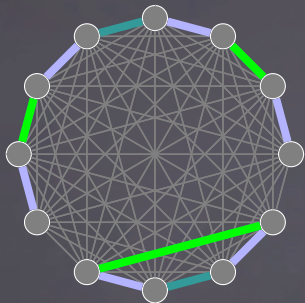


$$\sigma_p = 1, \sigma_c = 1?$$

### Approximate Exact Scaffolding

1. compute max-weight perfect matching  $S$   
 $\leadsto S \cup M$  is collection of cycles
2. "fix" all but lightest edge per cycle
3. repeatedly flip any lightest non-fix 4-cycle intersecting 2 cycles until at most  $\sigma_c + \sigma_p$  cycles remain
4. repeatedly remove lightest non-fix cycle-edge until at most  $\sigma_c$  cycles remain

## 2-Approximation in Dense Graphs



$$\sigma_p = 1, \sigma_c = 1?$$

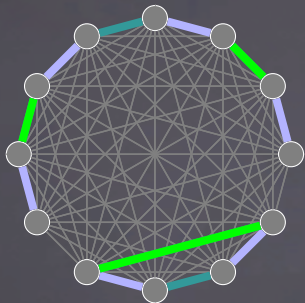
### Proof

Result  $S^*$  is a valid solution ✓

### Approximate Exact Scaffolding

1. compute max-weight perfect matching  $S$   
 $\leadsto S \cup M$  is collection of cycles
2. "fix" all but lightest edge per cycle
3. repeatedly flip any lightest non-fix 4-cycle intersecting 2 cycles until at most  $\sigma_c + \sigma_p$  cycles remain
4. repeatedly remove lightest non-fix cycle-edge until at most  $\sigma_c$  cycles remain

## 2-Approximation in Dense Graphs



$$\sigma_p = 1, \sigma_c = 1?$$

### Proof

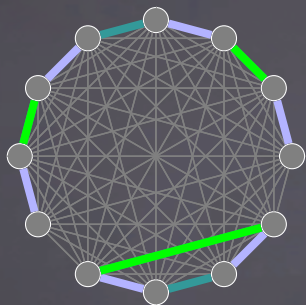
Result  $S^*$  is a valid solution ✓

$$\omega(S^*) \geq \omega(\text{fix}) \geq \omega(S)/2 \geq \text{OPT}/2$$

### Approximate Exact Scaffolding

1. compute max-weight perfect matching  $S$   
 $\leadsto S \cup M$  is collection of cycles
2. "fix" all but lightest edge per cycle
3. repeatedly flip any lightest non-fix 4-cycle intersecting 2 cycles until at most  $\sigma_c + \sigma_p$  cycles remain
4. repeatedly remove lightest non-fix cycle-edge until at most  $\sigma_c$  cycles remain

## 2-Approximation in Dense Graphs



$$\sigma_p = 1, \sigma_c = 1?$$

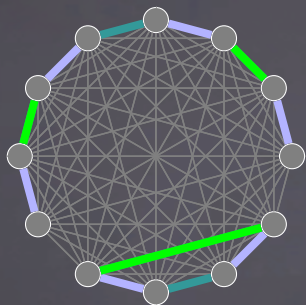
### Theorem

Exact Scaffolding in complete graphs can be 2-approximated in  $O(|V|^2)$  time.

### Approximate Exact Scaffolding

1. compute max-weight perfect matching  $S$   
 $\leadsto S \cup M$  is collection of cycles
2. "fix" all but lightest edge per cycle
3. repeatedly flip any lightest non-fix 4-cycle intersecting 2 cycles until at most  $\sigma_c + \sigma_p$  cycles remain
4. repeatedly remove lightest non-fix cycle-edge until at most  $\sigma_c$  cycles remain

## 2-Approximation in Dense Graphs



$$\sigma_p = 1, \sigma_c = 1?$$

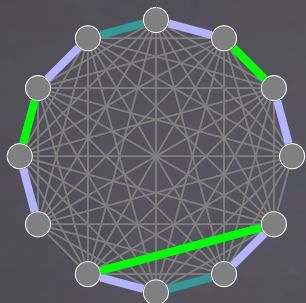
### Theorem

Exact Scaffolding in complete (Bipartite) graphs can be 2-approximated in  $O(|V|^2)$  time.

### Approximate Exact Scaffolding

1. compute max-weight perfect matching  $S$   
 $\leadsto S \cup M$  is collection of cycles
2. "fix" all but lightest edge per cycle
3. repeatedly flip any lightest non-fix 4-cycle intersecting 2 cycles until at most  $\sigma_c + \sigma_p$  cycles remain
4. repeatedly remove lightest non-fix cycle-edge until at most  $\sigma_c$  cycles remain

## 2-Approximation in Dense Graphs



$$\sigma_p = 1, \sigma_c = 1?$$

### Approximate Exact Scaffolding

1. compute max-weight perfect matching  $S$   
 $\leadsto S \cup M$  is collection of cycles
2. "fix" all but lightest edge per cycle
3. repeatedly flip any lightest non-fix 4-cycle intersecting 2 cycles until at most  $\sigma_c + \sigma_p$  cycles remain
4. repeatedly remove lightest non-fix cycle-edge until at most  $\sigma_c$  cycles remain

### Theorem

Exact Scaffolding in complete (Bipartite) graphs can be 2-approximated in  $O(|V|^2)$  time.

### Remark

For Scaffolding, replace Step 3 by either merging cycles or removing lightest edge, whatever loses less weight

# Exact Algorithms I: Brute Force

## Observation

An ordering of  $V(G)$  certifies YES-instances of Scaffolding.

# Exact Algorithms I: Brute Force



$[p, c, j]_i := \max_{\text{paths/cycles plus one path starting at } v_j} \text{weight collectible before } v_i \text{ with } p \notin c$

## Observation

An ordering of  $V(G)$  certifies YES-instances of Scaffolding.

# Exact Algorithms I: Brute Force



$[p, c, j]_i := \max_{\text{paths/cycles plus one path starting at } v_j} \text{weight collectible before } v_i \text{ with } p \notin c$

$$[p, c, j]_i = [p, c, j]_{i-2} + \omega(v_{i-2}v_{i-1}) \quad \text{if } j < i-2 \nmid v_{i-2}v_{i-1} \in E$$

## Observation

An ordering of  $V(G)$  certifies YES-instances of *Scaffolding*.

# Exact Algorithms I: Brute Force



$[p, c, j]_i := \max$  weight collectible before  $v_i$  with  $p \notin c$   
 paths/cycles plus one path starting at  $v_j$

$$[p, c, j]_i = [p, c, j]_{i-2} + \omega(v_{i-2}v_{i-1}) \quad \text{if } j < i-2 \nmid v_{i-2}v_{i-1} \in E$$

$$[p, c, i-1]_i = \max_{\substack{j < i-2 \\ j \text{ even}}} \begin{cases} [p-1, c, j]_{i-2} \end{cases}$$

## Observation

An ordering of  $V(G)$  certifies YES-instances of Scaffolding.

# Exact Algorithms I: Brute Force



$[p, c, j]_i := \max$  weight collectible before  $v_i$  with  $p \notin c$   
 paths/cycles plus one path starting at  $v_j$

$$[p, c, j]_i = [p, c, j]_{i-2} + \omega(v_{i-2}v_{i-1}) \quad \text{if } j < i-2 \nmid v_{i-2}v_{i-1} \in E$$

$$[p, c, i-1]_i = \max_{\substack{j < i-2 \\ j \text{ even}}} \begin{cases} [p-1, c, j]_{i-2} \\ [p, c-1, j]_{i-2} + \omega(v_j v_{i-2}) \quad \text{if } v_j v_{i-2} \in E \end{cases}$$

## Observation

An ordering of  $V(G)$  certifies YES-instances of Scaffolding.

# Exact Algorithms I: Brute Force



$[p, c, j]_i := \max$  weight collectible before  $v_i$  with  $p \notin c$   
 paths/cycles plus one path starting at  $v_j$

$$[p, c, j]_i = [p, c, j]_{i-2} + \omega(v_{i-2}v_{i-1}) \quad \text{if } j < i-2 \nmid v_{i-2}v_{i-1} \in E$$

$$[p, c, i-1]_i = \max_{\substack{j < i-2 \\ j \text{ even}}} \begin{cases} [p-1, c, j]_{i-2} \\ [p, c-1, j]_{i-2} + \omega(v_j v_{i-2}) \quad \text{if } v_j v_{i-2} \in E \end{cases}$$

## Observation

An ordering of  $V(G)$  certifies YES-instances of Scaffolding.

$\rightsquigarrow$  try all  $O(n!)$  certificates

# Exact Algorithms I: Brute Force



$[p, c, j]_i := \max$  weight collectible before  $v_i$  with  $p \notin c$   
paths/cycles plus one path starting at  $v_j$

$$[p, c, j]_i = [p, c, j]_{i-2} + \omega(v_{i-2}v_{i-1}) \quad \text{if } j < i-2 \nmid v_{i-2}v_{i-1} \in E$$

$$[p, c, i-1]_i = \max_{\substack{j < i-2 \\ j \text{ even}}} \begin{cases} [p-1, c, j]_{i-2} \\ [p, c-1, j]_{i-2} + \omega(v_j v_{i-2}) \quad \text{if } v_j v_{i-2} \in E \end{cases}$$

## Observation

An ordering of  $V(G)$  certifies YES-instances of Scaffolding.

$\rightsquigarrow$  try all  $O(n!)$  certificates

contigs force every other vertex  $\rightsquigarrow O(n!!)$

# Exact Algorithms I: Brute Force



$[p, c, j]_i := \max$  weight collectible before  $v_i$  with  $p \notin c$   
paths/cycles plus one path starting at  $v_j$

$$[p, c, j]_i = [p, c, j]_{i-2} + \omega(v_{i-2}v_{i-1}) \quad \text{if } j < i-2 \nmid v_{i-2}v_{i-1} \in E$$

$$[p, c, i-1]_i = \max_{\substack{j < i-2 \\ j \text{ even}}} \begin{cases} [p-1, c, j]_{i-2} \\ [p, c-1, j]_{i-2} + \omega(v_j v_{i-2}) \quad \text{if } v_j v_{i-2} \in E \end{cases}$$

## Observation

An ordering of  $V(G)$  certifies YES-instances of Scaffolding.

$\rightsquigarrow$  try all  $O(n!)$  certificates

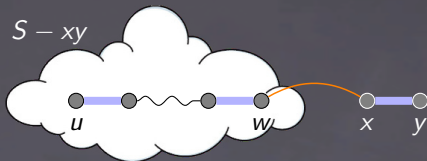
contigs force every other vertex  $\rightsquigarrow O(\sqrt{2}^n \cdot n/2!)$

# Exact Algorithms II: Dynamic Programming

## Semantics

$[S, p, c, u, v]$  = max weight collectible in  $G[S]$  by  $p$  alt. paths,  $c$  alt. cycles and an alt. path starting at  $u$  & ending at  $v$

# Exact Algorithms II: Dynamic Programming



## Semantics

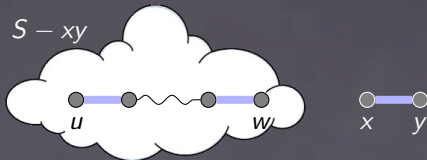
$[S, p, c, u, v] = \text{max weight collectible in } G[S] \text{ by } p \text{ alt. paths, } c \text{ alt. cycles and an alt. path starting at } u \text{ \& ending at } v$

## Computation

Let  $xy \in \mathcal{M}$ . Then,  $[\{xy\}, 0, 0, x, y] := 0$  and

$$[S, p, c, u, y] := \max_{\substack{w \in G[S-xy] \\ u \neq w}} [S - xy, p, c, u, w] + \omega(wx)$$

# Exact Algorithms II: Dynamic Programming



## Semantics

$[S, p, c, u, v] = \max$  weight collectible in  $G[S]$  by  $p$  alt. paths,  $c$  alt. cycles and an alt. path starting at  $u$  and ending at  $v$

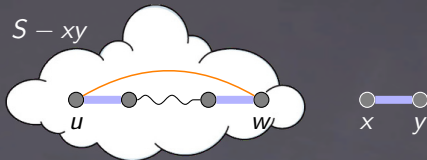
## Computation

Let  $xy \in \mathcal{M}$ . Then,  $[\{xy\}, 0, 0, x, y] := 0$  and

$$[S, p, c, u, y] := \max_{\substack{w \in G[S-xy] \\ u \neq w}} [S - xy, p, c, u, w] + \omega(wx)$$

$$[S, p, c, x, y] := \max_{u, w \in G[S-xy]} \begin{cases} [S - xy, p - 1, c, u, w] \\ \end{cases}$$

# Exact Algorithms II: Dynamic Programming



## Semantics

$[S, p, c, u, v] = \text{max weight collectible in } G[S] \text{ by } p \text{ alt. paths, } c \text{ alt. cycles and an alt. path starting at } u \text{ \& ending at } v$

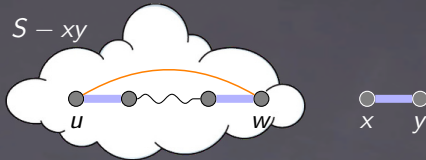
## Computation

Let  $xy \in \mathcal{M}$ . Then,  $[\{xy\}, 0, 0, x, y] := 0$  and

$$[S, p, c, u, y] := \max_{\substack{w \in G[S-xy] \\ u \neq w}} [S - xy, p, c, u, w] + \omega(wx)$$

$$[S, p, c, x, y] := \max_{u, w \in G[S-xy]} \begin{cases} [S - xy, p - 1, c, u, w] \\ [S - xy, p, c - 1, u, w] + \omega(wu) \end{cases} \text{ if } wu \in E(G) \setminus \mathcal{M}$$

# Exact Algorithms II: Dynamic Programming



## Semantics

$[S, p, c, u, v] = \text{max weight collectible in } G[S] \text{ by } p \text{ alt. paths, } c \text{ alt. cycles and an alt. path starting at } u \text{ \& ending at } v$

## Theorem

Scaffolding can be solved in  $O(\sqrt{2}^n n^3 \sigma_p \sigma_c)$  time.

# Sparse Graphs: Quasi-Forest

## Recall

- Scaffolding is hard in any sufficiently dense graph class
- Scaffolding is easy in trees

# Sparse Graphs: Quasi-Forest

## Recall

- Scaffolding is hard in any sufficiently dense graph class
- Scaffolding is easy in trees

## A Shot at Sparsity

$G$  is Quasi-forest  $\Leftrightarrow G - \mathcal{M}$  is forest

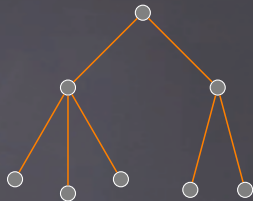
# Sparse Graphs: Quasi-Forest

## Recall

- Scaffolding is hard in any sufficiently dense graph class
- Scaffolding is easy in trees

## A Shot at Sparsity

$G$  is Quasi-forest  $\Leftrightarrow G - \mathcal{M}$  is forest



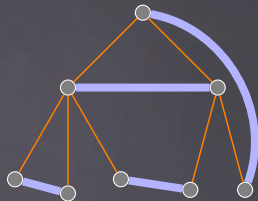
# Sparse Graphs: Quasi-Forest

## Recall

- Scaffolding is hard in any sufficiently dense graph class
- Scaffolding is easy in trees

## A Shot at Sparsity

$G$  is Quasi-forest  $\Leftrightarrow G - \mathcal{M}$  is forest



# Sparse Graphs: Quasi-Forest

## Recall

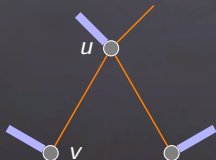
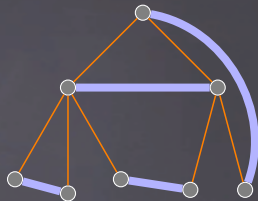
- Scaffolding is hard in any sufficiently dense graph class
- Scaffolding is easy in trees

## A Shot at Sparsity

$G$  is Quasi-forest  $\Leftrightarrow G - \mathcal{M}$  is forest

## Observation

Each leaf  $v$  of  $G - \mathcal{M}$  has degree 2 in  $G$   
 $\leadsto$  if unweighted, can we take both?



# Sparse Graphs: Quasi-Forest

## Recall

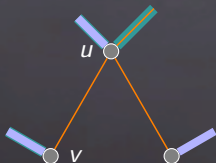
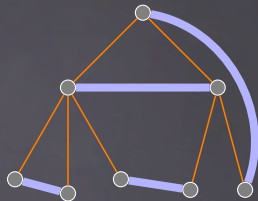
- Scaffolding is hard in any sufficiently dense graph class
- Scaffolding is easy in trees

## A Shot at Sparsity

$G$  is Quasi-forest  $\Leftrightarrow G - \mathcal{M}$  is forest

## Observation

Each leaf  $v$  of  $G - \mathcal{M}$  has degree 2 in  $G$   
 $\leadsto$  if unweighted, can we take both?



# Sparse Graphs: Quasi-Forest

## Recall

- Scaffolding is hard in any sufficiently dense graph class
- Scaffolding is easy in trees

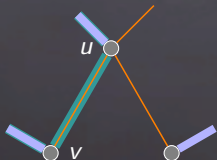
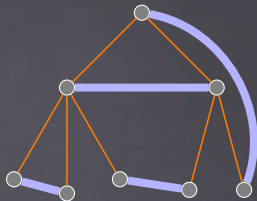
## A Shot at Sparsity

$G$  is Quasi-forest  $\Leftrightarrow G - \mathcal{M}$  is forest

## Observation

Each leaf  $v$  of  $G - \mathcal{M}$  has degree 2 in  $G$

$\rightsquigarrow$  if unweighted, can we take both?



## Observation

- $v$  in path  $\nsubseteq u$  in cycle  $\rightsquigarrow$  1 path ✓
- $v$  in path  $\nsubseteq u$  in path  $\rightsquigarrow$  2 paths ✓

# Sparse Graphs: Quasi-Forest

## Recall

- Scaffolding is hard in any sufficiently dense graph class
- Scaffolding is easy in trees

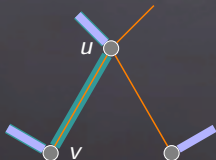
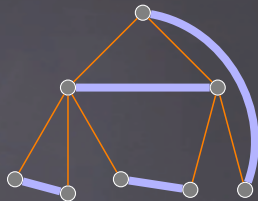
## A Shot at Sparsity

$G$  is Quasi-forest  $\Leftrightarrow G - \mathcal{M}$  is forest

## Observation

Each leaf  $v$  of  $G - \mathcal{M}$  has degree 2 in  $G$

$\rightsquigarrow$  if unweighted, can we take both?



## Observation

- $v$  in path  $\nsubseteq u$  in cycle  $\rightsquigarrow$  1 path ✓
- $v$  in path  $\nsubseteq u$  in path  $\rightsquigarrow$  2 paths ✓  
unless it's the same path!

# Sparse Graphs: Quasi-Forest

## Recall

- Scaffolding is hard in any sufficiently dense graph class
- Scaffolding is easy in trees

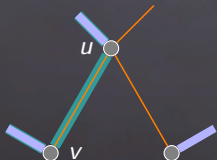
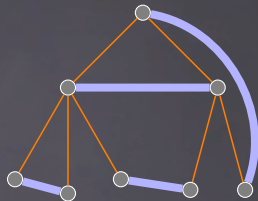
## A Shot at Sparsity

$G$  is Quasi-forest  $\Leftrightarrow G - \mathcal{M}$  is forest

## Observation

Each leaf  $v$  of  $G - \mathcal{M}$  has degree 2 in  $G$

$\rightsquigarrow$  if unweighted, can we take both? ✗



## Observation

- $v$  in path  $\nsubseteq u$  in cycle  $\rightsquigarrow$  1 path ✓
- $v$  in path  $\nsubseteq u$  in path  $\rightsquigarrow$  2 paths ✓

# Sparse Graphs: Quasi-Forest

## Recall

- Scaffolding is hard in any sufficiently dense graph class
- Scaffolding is easy in trees

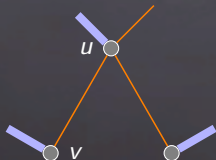
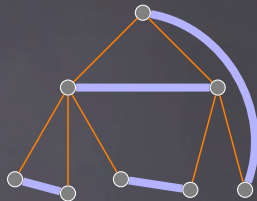
## A Shot at Sparsity

$G$  is Quasi-forest  $\Leftrightarrow G - \mathcal{M}$  is forest

## Observation

Each leaf  $v$  of  $G - \mathcal{M}$  has degree 2 in  $G$

$\leadsto$  if  $\sigma_p = 0$ , we have to take both!



# Sparse Graphs: Quasi-Forest

## Recall

- Scaffolding is hard in any sufficiently dense graph class
- Scaffolding is easy in trees

## A Shot at Sparsity

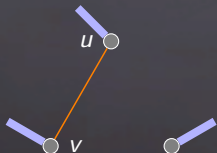
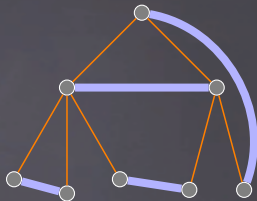
$G$  is Quasi-forest  $\Leftrightarrow G - \mathcal{M}$  is forest

## Observation

Each leaf  $v$  of  $G - \mathcal{M}$  has degree 2 in  $G$

$\rightsquigarrow$  if  $\sigma_p = 0$ , we have to take both!

$\rightsquigarrow$  remove all non-matching edges from parent  $u$ , except  $uv$



# Sparse Graphs: Quasi-Forest

## Recall

- Scaffolding is hard in any sufficiently dense graph class
- Scaffolding is easy in trees

## A Shot at Sparsity

$G$  is Quasi-forest  $\Leftrightarrow G - \mathcal{M}$  is forest

## Observation

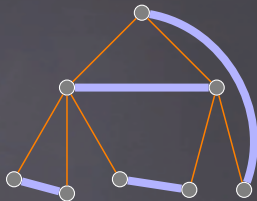
Each leaf  $v$  of  $G - \mathcal{M}$  has degree 2 in  $G$

$\leadsto$  if  $\sigma_p = 0$ , we have to take both!

$\leadsto$  remove all non-matching edges from parent  $u$ , except  $uv$

## Corollary

Scaffolding can be solved in  $O(n)$  on quasi-forests if  $\sigma_p = 0$ .



# Sparse Graphs: Quasi-Forest

## Recall

- Scaffolding is hard in any sufficiently dense graph class
- Scaffolding is easy in trees

## A Shot at Sparsity

$G$  is Quasi-forest  $\Leftrightarrow G - \mathcal{M}$  is forest

## Observation

Each leaf  $v$  of  $G - \mathcal{M}$  has degree 2 in  $G$

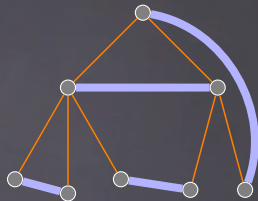
$\rightsquigarrow$  if  $\sigma_p = 0$ , we have to take both!

$\rightsquigarrow$  remove all non-matching edges from parent  $u$ , except  $uv$

## Corollary

Scaffolding can be solved in  $O(n)$  on quasi-forests if  $\sigma_p = 0$ .

Scaffolding can be solved in  $O(n^{2\sigma_p+1})$  in quasi-forests.



# Sparse Graphs: Quasi-Forest

## Recall

- Scaffolding is hard in any sufficiently dense graph class
- Scaffolding is easy in trees

## A Shot at Sparsity

$G$  is Quasi-forest  $\Leftrightarrow G - \mathcal{M}$  is forest

## Observation

Each leaf  $v$  of  $G - \mathcal{M}$  has degree 2 in  $G$

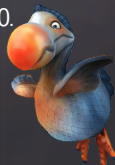
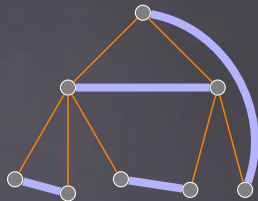
$\rightsquigarrow$  if  $\sigma_p = 0$ , we have to take both!

$\rightsquigarrow$  remove all non-matching edges from parent  $u$ , except  $uv$

## Corollary

Scaffolding can be solved in  $O(n)$  on quasi-forests if  $\sigma_p = 0$ .

Scaffolding can be solved in  $O(n^{2\sigma_p+1})$  in quasi-forests.



But is it even NP-hard?

# Sparse Graphs: Quasi-Forest

## Weighted 2-SAT

**Input:**  $\varphi$  on  $X$  in 2-CNF, weights  $w : X \times \{0, 1\} \rightarrow \mathbb{N}$ ,  $k \in \mathbb{N}$

**Question:** is there a satisfying assignment for  $\varphi$  of weight  $\leq k$ ?

# Sparse Graphs: Quasi-Forest

## Weighted 2-SAT

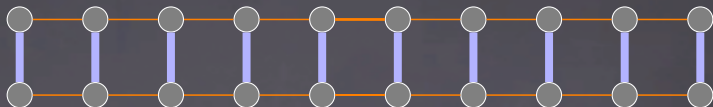
Input:  $\varphi$  on  $X$  in 2-CNF, weights  $w : X \times \{0, 1\} \rightarrow \mathbb{N}$ ,  $k \in \mathbb{N}$

Question: is there a satisfying assignment for  $\varphi$  of weight  $\leq k$ ?

## Remark

Independent Set is special case of Weighted 2-SAT

## Sparse Graphs: Quasi-Forest



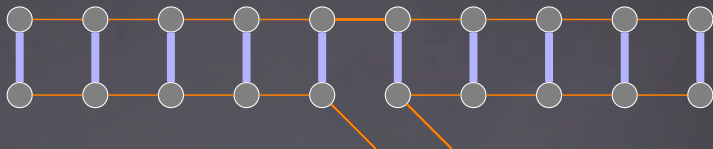
# Sparse Graphs: Quasi-Forest



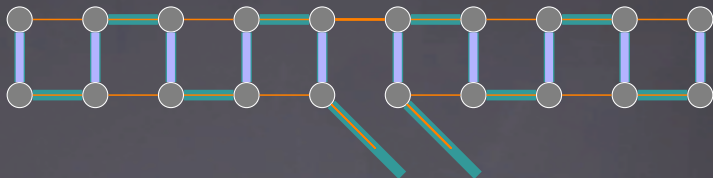
# Sparse Graphs: Quasi-Forest



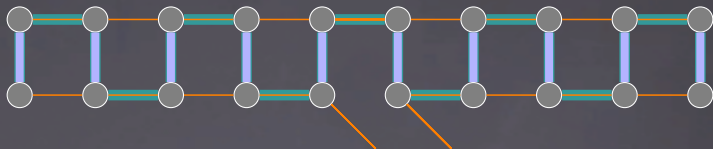
## Sparse Graphs: Quasi-Forest



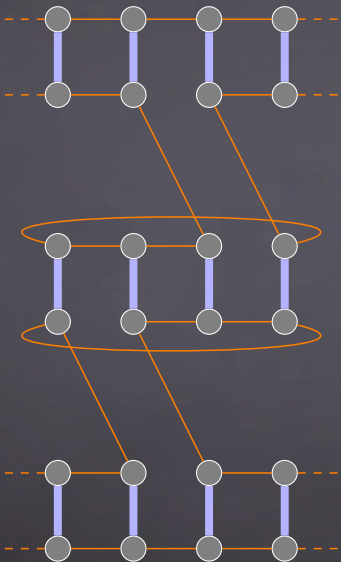
# Sparse Graphs: Quasi-Forest



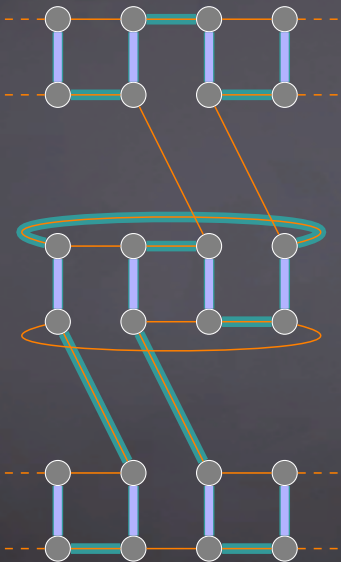
# Sparse Graphs: Quasi-Forest



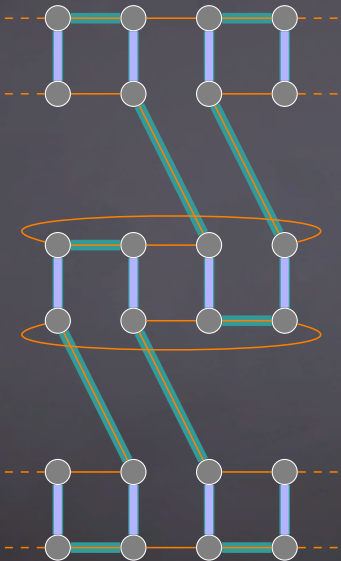
# Sparse Graphs: Quasi-Forest



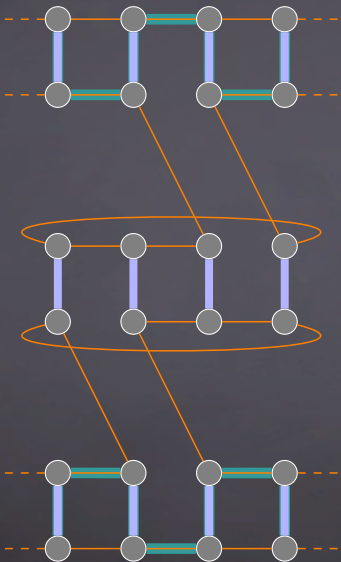
# Sparse Graphs: Quasi-Forest



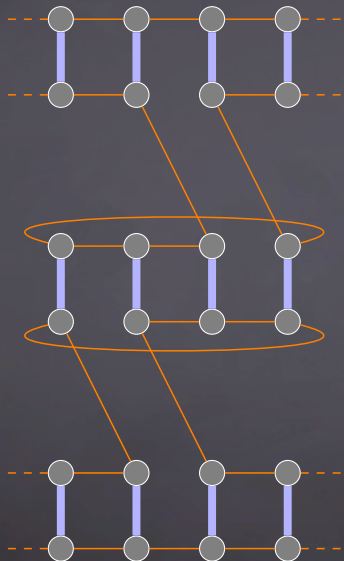
# Sparse Graphs: Quasi-Forest



# Sparse Graphs: Quasi-Forest



# Sparse Graphs: Quasi-Forest



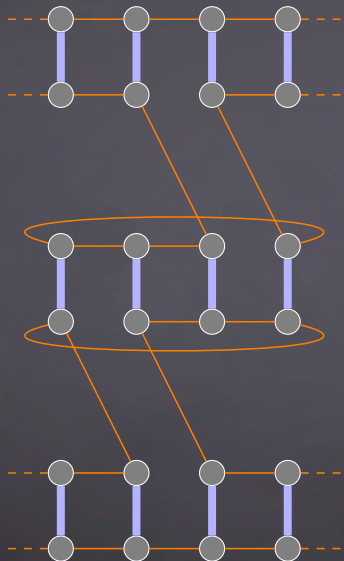
## Observation

$\exists$  weight- $k$  satisfying assignment

$\Leftrightarrow$

$\exists$  weight- $k$  cover with  $\leq n$   
alternating paths

# Sparse Graphs: Quasi-Forest



## Observation

$\exists$  weight- $k$  satisfying assignment

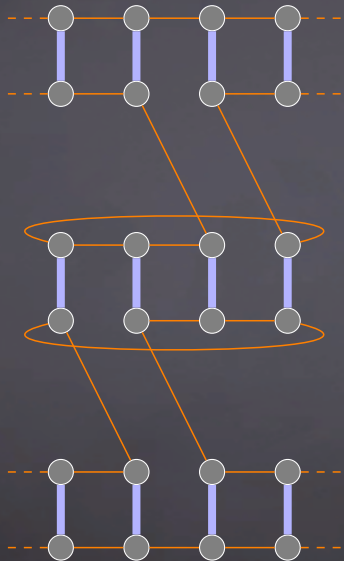
$\Leftrightarrow$

$\exists$  weight- $k$  cover with  $\leq n$   
alternating paths

## Theorem

Scaffolding is NP-hard even if  $G - \mathcal{M}$   
is a collection of paths with weights  
0/1

# Sparse Graphs: Quasi-Forest



## Observation

$\exists$  weight- $k$  satisfying assignment

$\Leftrightarrow$

$\exists$  weight- $k$  cover with  $\leq n$   
alternating paths

## Theorem

Scaffolding is NP-hard even if  $G - \mathcal{M}$   
is a collection of paths with weights  
0/1

## Corollary

no  $2^{o(n+m)}$ -time algorithm (ETH)

no  $n^{o(k)}$ -time algorithm ( $\text{FPT} \neq \text{W}[t]$ )

# Other Forms of Tree-Likeness

## Tree Decompositions

tree  $T$ , each vertex  $i$  associated to some  $X_i \subseteq V(G)$  s.t.

1.  $\forall e \in E(G)$ , there is some  $i \in V(T)$  with  $e \in X_i$
2.  $\forall v \in V(G)$ , bags containing  $v$  induce a connected subtree

**treewidth**  $tw = \text{size of largest bag} - 1$

# Other Forms of Tree-Likeness

## Tree Decompositions

tree  $T$ , each vertex  $i$  associated to some  $X_i \subseteq V(G)$  s.t.

1.  $\forall e \in E(G)$ , there is some  $i \in V(T)$  with  $e \in X_i$
  2.  $\forall v \in V(G)$ , bags containing  $v$  induce a connected subtree
- treewidth**  $tw = \text{size of largest bag} - 1$

## Hope

Practical instances of **Scaffolding** have **low treewidth** (they originate from linear structure)

# Other Forms of Tree-Likeness

## Tree Decompositions

tree  $T$ , each vertex  $i$  associated to some  $X_i \subseteq V(G)$  s.t.

1.  $\forall e \in E(G)$ , there is some  $i \in V(T)$  with  $e \in X_i$
  2.  $\forall v \in V(G)$ , bags containing  $v$  induce a connected subtree
- treewidth**  $tw = \text{size of largest bag} - 1$

## Hope

Practical instances of **Scaffolding** have **low treewidth** (they originate from linear structure)

## Nice Decompositions

**Leaf**:  $X = \emptyset$

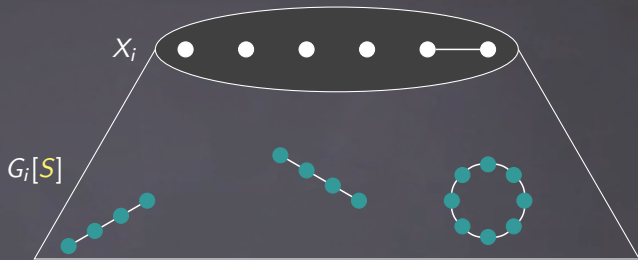
**Introduce  $v$** :  $i$  has single child  $j$  and  $X_i \setminus X_j = \{v\}$

**Forget  $v$** :  $i$  has single child  $j$  and  $X_j \setminus X_i = \{v\}$

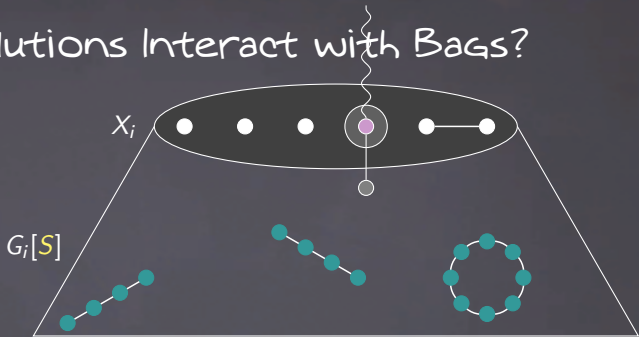
**Introduce  $uv$** :  $i$  has single child  $j$  and  $uv \subseteq X_i = X_j$   
(each edge introduced exactly once)

**Join**:  $i$  has 2 children  $j$  and  $\ell$  and  $X_i = X_j = X_\ell$

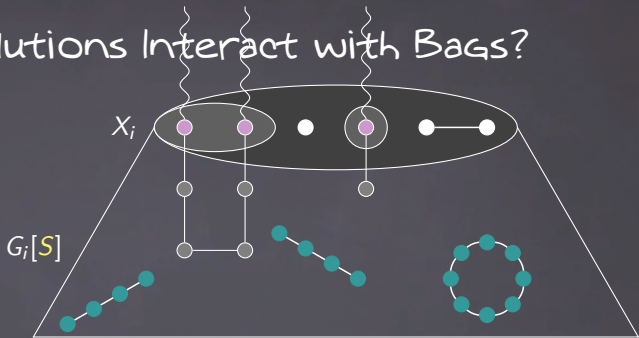
# How do Solutions Interact with Bags?



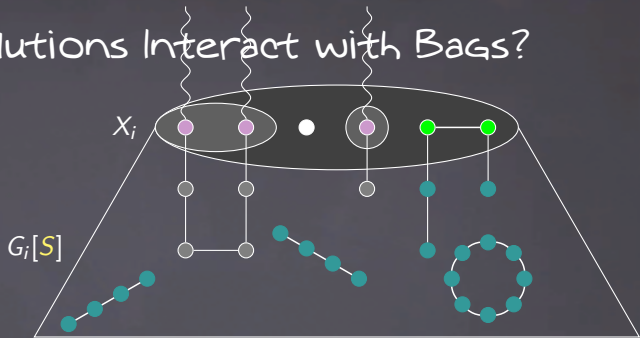
# How do Solutions Interact with Bags?



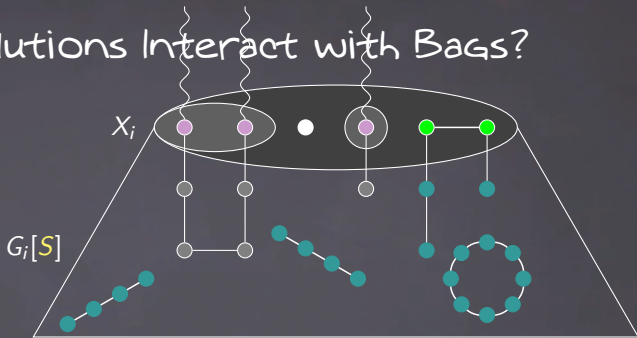
# How do Solutions Interact with Bags?



# How do Solutions Interact with Bags?



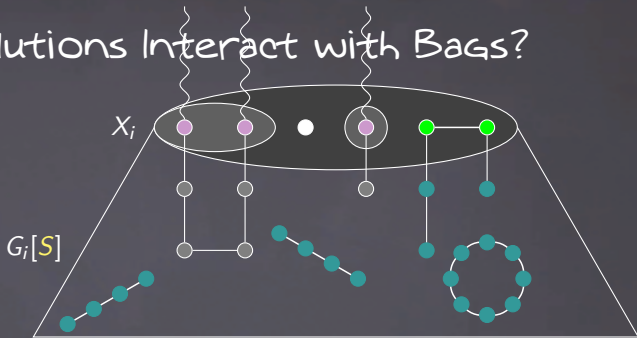
# How do Solutions Interact with Bags?



## Ingredients

- degree-function  $d : X \rightarrow \{0, 1, 2\}$
- "pairing"  $\subseteq \binom{X}{2} \cup X$
- $\#$ paths and  $\#$ cycles completed "Below the Bag"

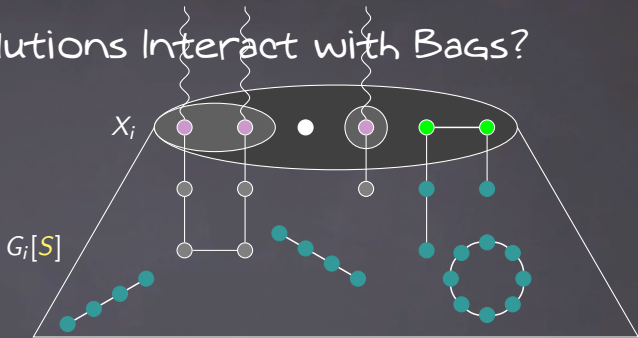
# How do Solutions Interact with Bags?



## Ingredients

- degree-function  $d : X \rightarrow \{0, 1, 2\}$
- "pairing"  $\subseteq \binom{X}{2} \cup X \rightsquigarrow \# \text{matchings possibilities} \rightsquigarrow O(|X|^{|X|/2})$
- $\# \text{paths}$  and  $\# \text{cycles}$  completed "Below the Bag"

# How do Solutions Interact with Bags?



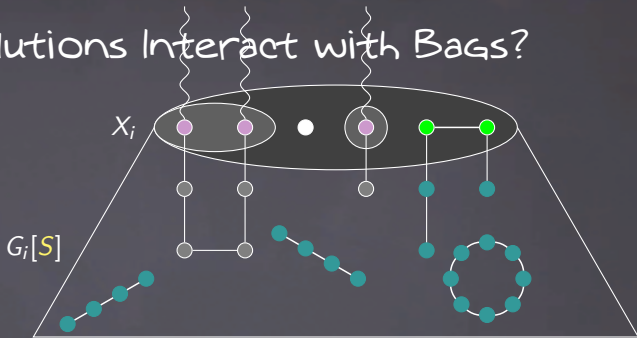
## Ingredients

- degree-function  $d : X \rightarrow \{0, 1, 2\}$
- "pairing"  $\subseteq \binom{X}{2} \cup X \rightsquigarrow \# \text{matchings possibilities} \rightsquigarrow O(|X|^{X/2})$
- $\# \text{paths}$  and  $\# \text{cycles}$  completed "Below the Bag"

## Semantics

- $[d, P, p, c]_i = \max$  weight of any  $S$  with  $M \cap E(G_i) \subseteq S \subseteq E(G_i)$  and
1. each vertex  $v \in X_i$  has degree  $d(v)$  in  $G_i[S]$ ,
  2. for each  $uv \in P$ ,  $G_i[S]$  contains an alternating path...
    - $u = v$ : ... from  $u$  avoiding  $d^{-1}(1)$
    - $u \neq v$ : ... from  $u$  to  $v$
  3.  $G_i[S]$  contains  $p$  alt. paths  $\nleftrightarrow c$  alt. cycles avoiding  $d^{-1}(1)$

# How do Solutions Interact with Bags?



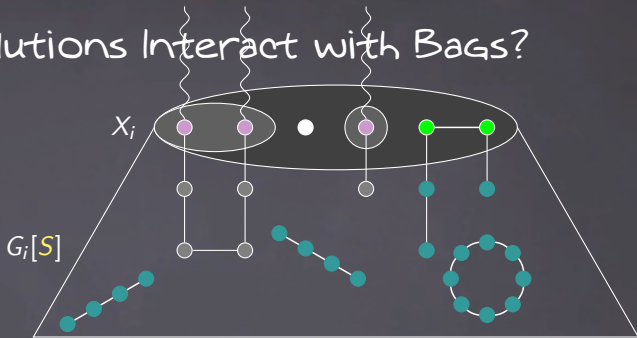
## Ingredients

- degree-function  $d : X \rightarrow \{0, 1, 2\}$
- "pairing"  $\subseteq \binom{X}{2} \cup X \rightsquigarrow \# \text{matchings possibilities} \rightsquigarrow O(|X|^{|X|/2})$
- $\# \text{paths}$  and  $\# \text{cycles}$  completed "below the bag"

## Leaf Bag

$$[\emptyset, \emptyset, 0, 0]_i = 0$$

# How do Solutions Interact with Bags?



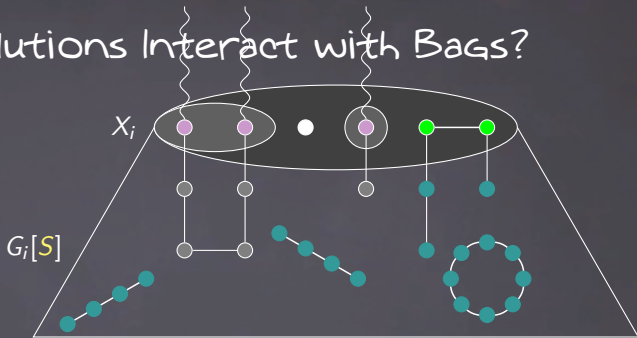
## Ingredients

- degree-function  $d : X \rightarrow \{0, 1, 2\}$
- "pairing"  $\subseteq \binom{X}{2} \cup X \rightsquigarrow \# \text{matchings possibilities} \rightsquigarrow O(|X|^{X/2})$
- $\# \text{paths}$  and  $\# \text{cycles}$  completed "Below the Bag"

## Introduce $v$ (single child $j$ )

$$[d, P, p, c]_i = [d|_{v \rightarrow \perp}, P, p, c]_j$$

# How do Solutions Interact with Bags?

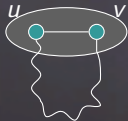


## Ingredients

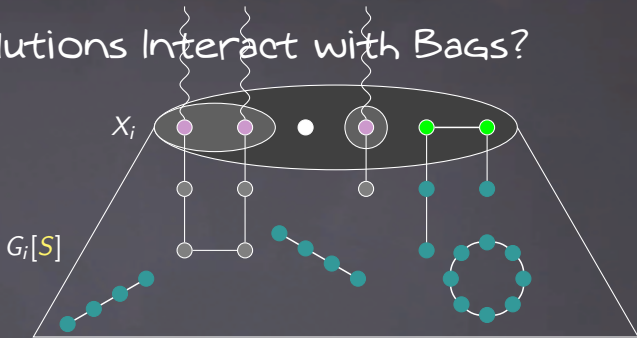
- degree-function  $d : X \rightarrow \{0, 1, 2\}$
- "pairing"  $\subseteq \binom{X}{2} \cup X \rightsquigarrow \# \text{matchings possibilities} \rightsquigarrow O(|X|^{X/2})$
- $\# \text{paths}$  and  $\# \text{cycles}$  completed "Below the Bag"

Introduce  $uv$  (single child  $j$ )

Case I:  $d(u) = d(v) = 2$



# How do Solutions Interact with Bags?

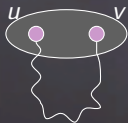


## Ingredients

- degree-function  $d : X \rightarrow \{0, 1, 2\}$
- "pairing"  $\subseteq \binom{X}{2} \cup X \rightsquigarrow \# \text{matchings possibilities} \rightsquigarrow O(|X|^{|X|/2})$
- $\# \text{paths}$  and  $\# \text{cycles}$  completed "Below the Bag"

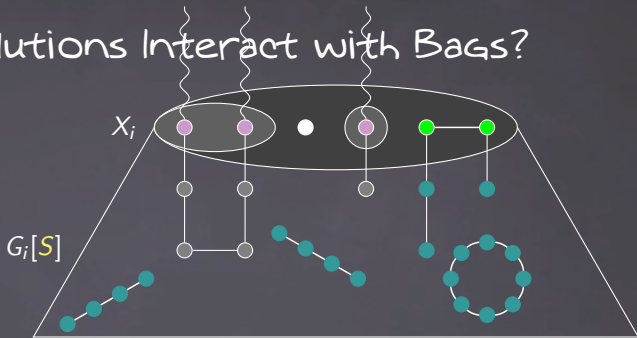
Introduce  $uv$  (single child  $j$ )

Case I:  $d(u) = d(v) = 2$



$$[d, P, p, c]_i = [d|_{u \rightarrow 1, v \rightarrow 1}, P + uv, p, c - 1]_j$$

# How do Solutions Interact with Bags?



## Ingredients

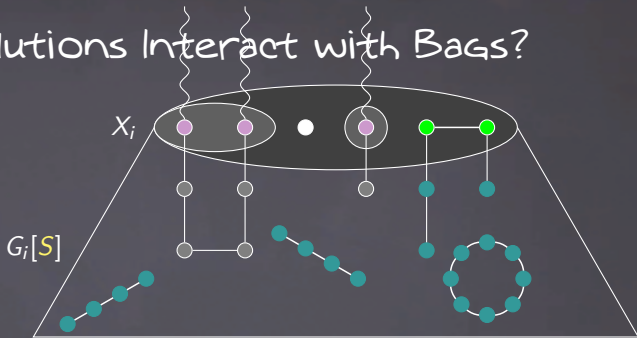
- degree-function  $d : X \rightarrow \{0, 1, 2\}$
- "pairing"  $\subseteq \binom{X}{2} \cup X \rightsquigarrow \# \text{matchings possibilities} \rightsquigarrow O(|X|^{|X|/2})$
- $\# \text{paths}$  and  $\# \text{cycles}$  completed "Below the Bag"

Introduce  $uv$  (single child  $j$ )

Case I:  $d(u) = d(v) = 2$



# How do Solutions Interact with Bags?



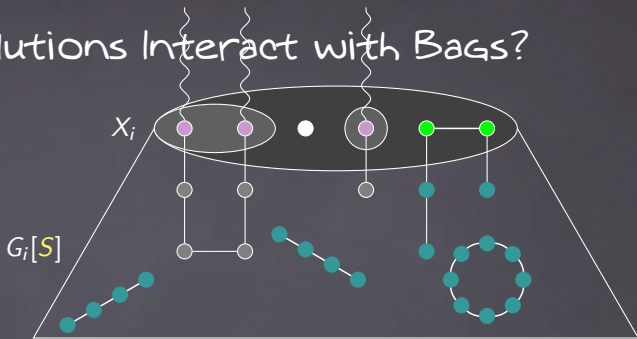
## Ingredients

- degree-function  $d : X \rightarrow \{0, 1, 2\}$
- "pairing"  $\subseteq \binom{X}{2} \cup X \rightsquigarrow \# \text{matchings possibilities} \rightsquigarrow O(|X|^{X/2})$
- $\# \text{paths}$  and  $\# \text{cycles}$  completed "Below the Bag"

Forget  $v$  (single child  $j$ )

$$[d, P, p, c]_i = \max \left\{ \right.$$

# How do Solutions Interact with Bags?



## Ingredients

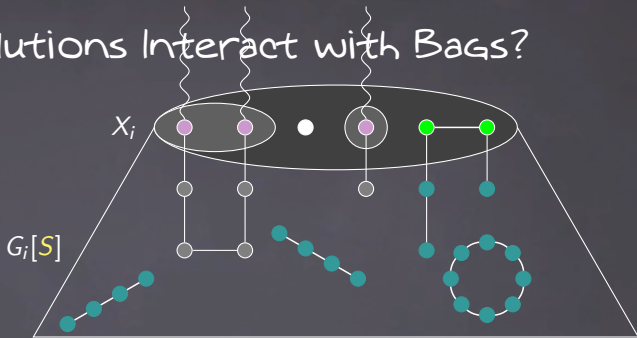
- degree-function  $d : X \rightarrow \{0, 1, 2\}$
- "pairing"  $\subseteq \binom{X}{2} \cup X \rightsquigarrow \# \text{matchings possibilities} \rightsquigarrow O(|X|^{|X|/2})$
- $\# \text{paths}$  and  $\# \text{cycles}$  completed "Below the Bag"

Forget  $v$  (single child  $j$ )



$$[d, P, p, c]_i = \max \left\{ \begin{array}{l} [d|_{v \rightarrow 1}, P + vv, p - 1, c]_j \end{array} \right.$$

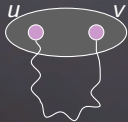
# How do Solutions Interact with Bags?



## Ingredients

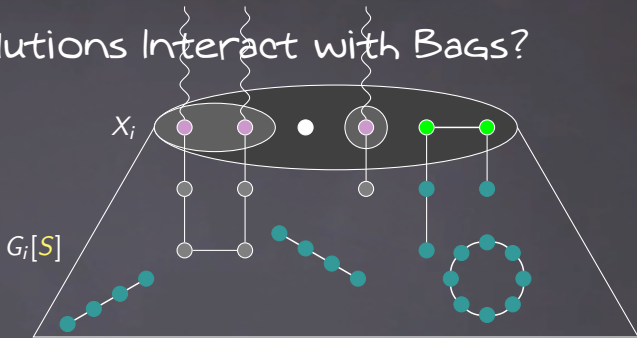
- degree-function  $d : X \rightarrow \{0, 1, 2\}$
- "pairing"  $\subseteq \binom{X}{2} \cup X \rightsquigarrow \# \text{matchings possibilities} \rightsquigarrow O(|X|^{|X|/2})$
- $\# \text{paths}$  and  $\# \text{cycles}$  completed "Below the Bag"

Forget  $v$  (single child  $j$ )



$$[d, P, p, c]_i = \max \begin{cases} [d|_{v \rightarrow 1}, P + vv, p - 1, c]_j \\ \max_{uu \in P} [d|_{v \rightarrow 1}, (P - uu) + uv, p, c]_j \end{cases}$$

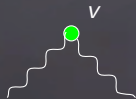
# How do Solutions Interact with Bags?



## Ingredients

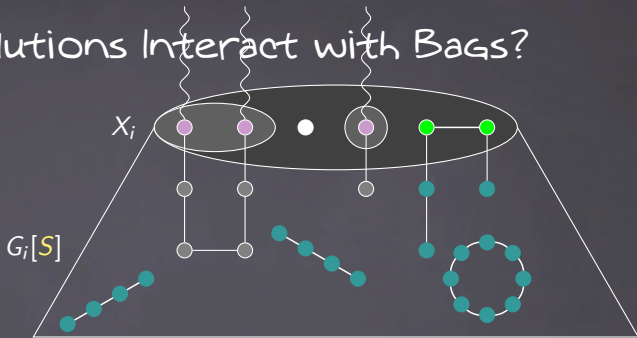
- degree-function  $d : X \rightarrow \{0, 1, 2\}$
- "pairing"  $\subseteq \binom{X}{2} \cup X \rightsquigarrow \# \text{matchings possibilities} \rightsquigarrow O(|X|^{|X|/2})$
- $\# \text{paths}$  and  $\# \text{cycles}$  completed "Below the Bag"

Forget  $v$  (single child  $j$ )



$$[d, P, p, c]_i = \max \begin{cases} [d|_{v \rightarrow 1}, P + vv, p - 1, c]_j \\ \max_{uu \in P} [d|_{v \rightarrow 1}, (P - uu) + uv, p, c]_j \\ \max_{x \in \{0, 2\}} [d|_{v \rightarrow x}, P, p, c]_j \end{cases}$$

# How do Solutions Interact with Bags?



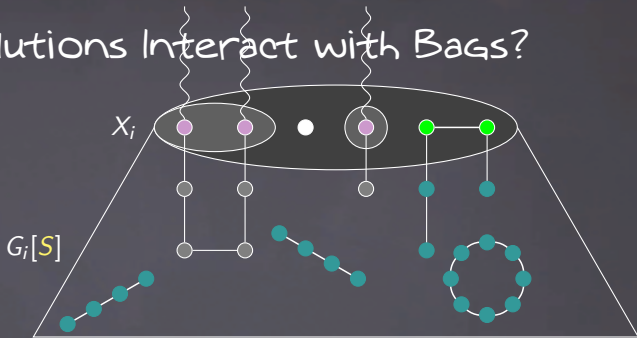
## Ingredients

- degree-function  $d : X \rightarrow \{0, 1, 2\}$
- "pairing"  $\subseteq \binom{X}{2} \cup X \rightsquigarrow \# \text{matchings possibilities} \rightsquigarrow O(|X|^{|X|/2})$
- $\# \text{paths}$  and  $\# \text{cycles}$  completed "Below the Bag"

## Join Bag (children $j \neq \ell$ )

$$[d, P, p, c]_i = \max_{d_j, P_j, p_j, c_j} \max_{\substack{P_\ell \\ P_j \sqcup P_\ell = P}} [d_j, P_j, p_j, c_j]_j + [d - d_j, P - P_j, p - p_j, c - c_j]_\ell$$

# How do Solutions Interact with Bags?



## Ingredients

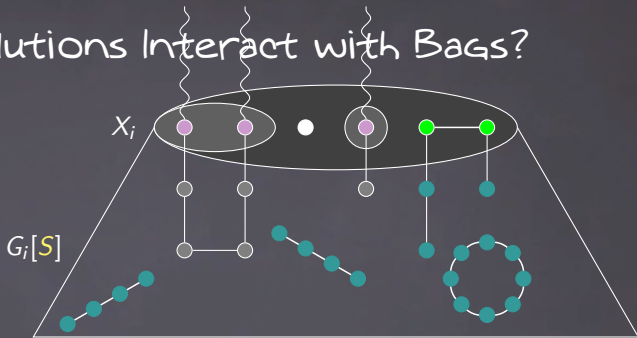
- degree-function  $d : X \rightarrow \{0, 1, 2\}$
- "pairing"  $\subseteq \binom{X}{2} \cup X \rightsquigarrow \# \text{matchings possibilities} \rightsquigarrow O(|X|^{|X|/2})$
- $\# \text{paths and } \# \text{cycles completed "Below the Bag"}$

## Join Bag (children $j \neq \ell$ )

$$[d, P, p, c]_i = \max_{d_j, P_j, p_j, c_j} \max_{\substack{P_\ell \\ P_j \sqcup P_\ell = P}} [d_j, P_j, p_j, c_j]_j + [d - d_j, P_\ell, p - p_j, c - c_j]_\ell$$

$$\rightsquigarrow O(3^{\text{tw}} \cdot \text{tw}^{\text{tw}/2} \cdot \sigma_p \cdot \sigma_c) \text{ table entries}$$

# How do Solutions Interact with Bags?



## Ingredients

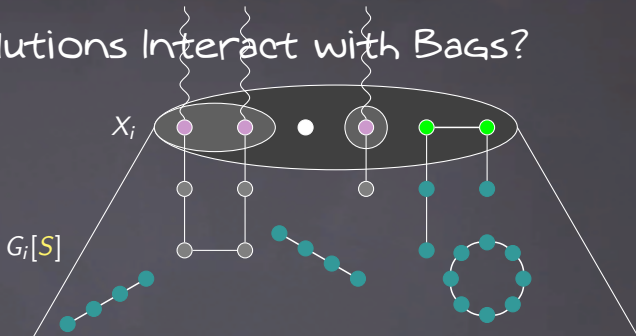
- degree-function  $d : X \rightarrow \{0, 1, 2\}$
- "pairing"  $\subseteq \binom{X}{2} \cup X \rightsquigarrow \# \text{matchings possibilities} \rightsquigarrow O(|X|^{|X|/2})$
- $\# \text{paths and } \# \text{cycles completed "Below the Bag"}$

## Join Bag (children $j \neq \ell$ )

$$[d, P, p, c]_i = \max_{d_j, P_j, p_j, c_j} \max_{\substack{P_\ell \\ P_j \sqcup P_\ell = P}} [d_j, P_j, p_j, c_j]_j + [d - d_j, P_\ell, p - p_j, c - c_j]_\ell$$

$\rightsquigarrow O(2^{\text{tw}} \cdot \text{tw}^{\text{tw}/2} \cdot \sigma_p \cdot \sigma_c)$  table entries

# How do Solutions Interact with Bags?



## Ingredients

- degree-function  $d : X \rightarrow \{0, 1, 2\}$
- "pairing"  $\subseteq \binom{X}{2} \cup X \rightsquigarrow \# \text{matchings possibilities} \rightsquigarrow O(|X|^{|X|/2})$
- $\# \text{paths and } \# \text{cycles completed "Below the Bag"}$

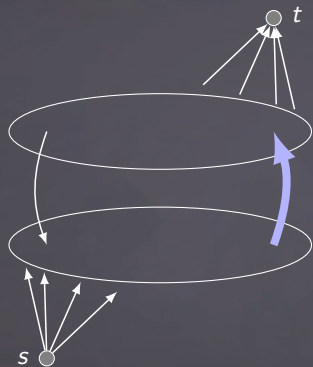
## Join Bag (children $j \neq \ell$ )

$$[d, P, p, c]_i = \max_{d_j, P_j, p_j, c_j} \max_{\substack{P_\ell \\ P_j \sqcup P_\ell = P}} [d_j, P_j, p_j, c_j]_j + [d - d_j, P_\ell, p - p_j, c - c_j]_\ell$$

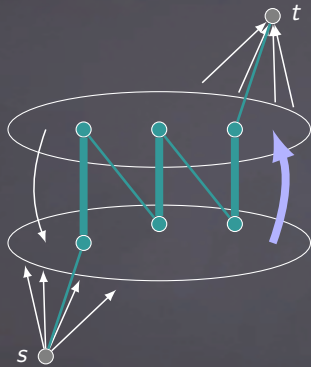
$\rightsquigarrow O(2^{\text{tw}} \cdot \text{tw}^{\text{tw}/2} \cdot \sigma_p \cdot \sigma_c)$  table entries and  $O((\text{tw} + 2)^{\text{tw}} \cdot \sigma_p \cdot \sigma_c \cdot n)$  time



# Integer Linear Program Formulation

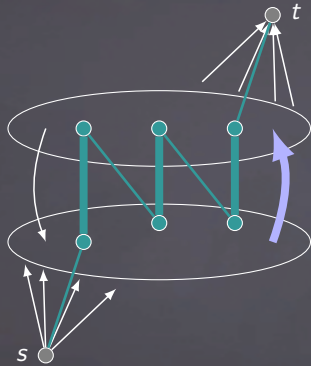


# Integer Linear Program Formulation



- chromosomes = disjoint  $s$ - $t$ -paths

# Integer Linear Program Formulation



- chromosomes = disjoint  $s$ - $t$ -paths

- Bin. variables

$$y_{uv} = 1 \Leftrightarrow u \rightarrow v \text{ used}$$

$$x_{\{u,v\}} = y_{uv} + y_{vu}$$

- force contigs:

$$\forall_{uv \in M} x_{uv} = 1$$

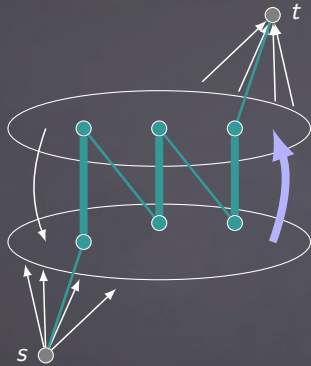
- path preservation:

$$\forall_{u \neq s, t} \sum_v y_{vu} = \sum_v y_{uv}$$

- path bounds:

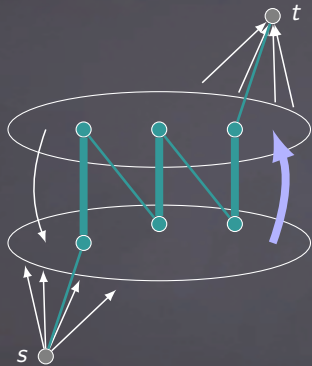
$$\sum_v y_{vt} \leq \sigma$$

# Integer Linear Program Formulation



- chromosomes = disjoint  $s$ - $t$ -paths
- Bin. variables  $y_{uv} = 1 \Leftrightarrow u \rightarrow v$  used  
 $x_{\{u,v\}} = y_{uv} + y_{vu}$
- force contigs:  $\forall_{uv \in M} x_{uv} = 1$
- path preservation:  $\forall_{u \neq s, t} \sum_v y_{vu} = \sum_v y_{uv}$
- path bounds:  $\sum_v y_{vt} \leq \sigma$
- forbid cycles (row generation via callback):  
 $\forall \text{ cycle } C: \sum_{uv \in C} y_{uv} < |C|$

# Integer Linear Program Formulation



- chromosomes = disjoint  $s$ - $t$ -paths

- Bin. variables

$$y_{uv} = 1 \Leftrightarrow u \rightarrow v \text{ used}$$

$$x_{\{u,v\}} = y_{uv} + y_{vu}$$

- force contigs:

$$\forall_{uv \in M} x_{uv} = 1$$

- path preservation:

$$\forall_{u \neq s, t} \sum_v y_{vu} = \sum_v y_{uv}$$

- path bounds:

$$\sum_v y_{vt} \leq \sigma$$

- forbid cycles (row generation via callback):

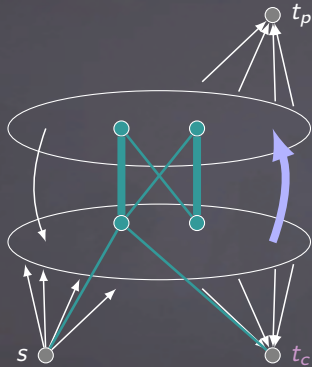
$\forall$  cycle  $C$ :

$$\sum_{uv \in C} y_{uv} < |C|$$

- Objective:

$$\max \sum_{e \in E} x_{\{u,v\}} \cdot \omega(e)$$

# Integer Linear Program Formulation



- chromosomes = disjoint  $s$ - $t$ -paths

- Bin. variables

$$y_{uv} = 1 \Leftrightarrow u \rightarrow v \text{ used}$$

$$x_{\{u,v\}} = y_{uv} + y_{vu}$$

- force contigs:

$$\forall_{uv \in M} x_{uv} = 1$$

- path preservation:

$$\forall_{u \neq s, t} \sum_v y_{vu} = \sum_v y_{uv}$$

- path bounds:

$$\sum_v y_{vt} \leq \sigma$$

- forbid cycles (row generation via callback):

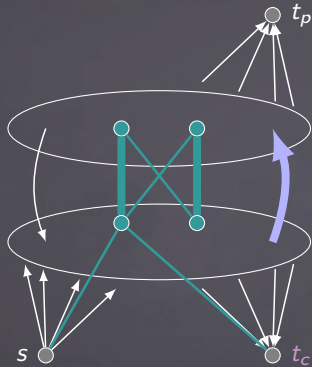
$\forall$  cycle  $C$ :

$$\sum_{uv \in C} y_{uv} < |C|$$

- Objective:

$$\max \sum_{e \in E} x_{\{u,v\}} \cdot \omega(e)$$

# Integer Linear Program Formulation



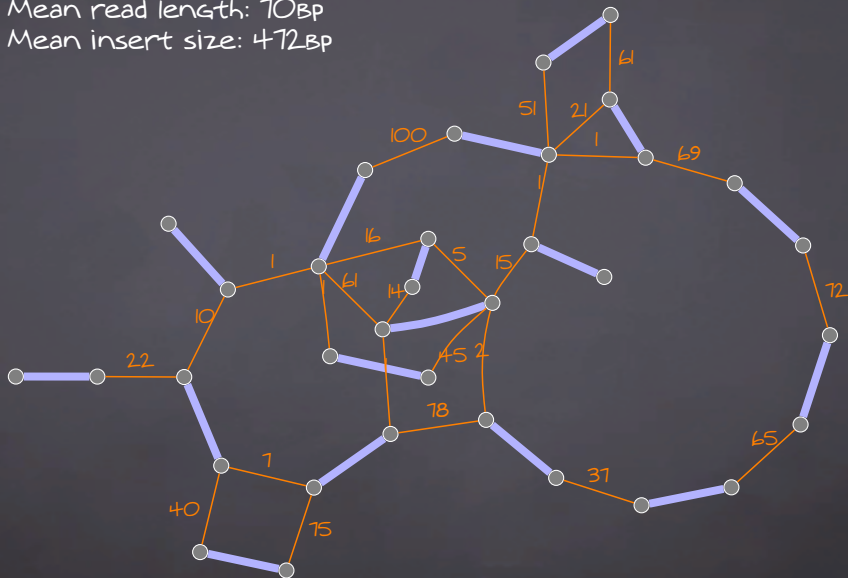
- chromosomes = disjoint  $s$ - $\{t_p, t_c\}$ -paths
- Bin. variables
 
$$y_{uv} = 1 \Leftrightarrow u \rightarrow v \text{ used}$$

$$x_{\{u,v\}} = y_{uv} + y_{vu}$$
- force contigs:
 
$$\forall_{uv \in M} x_{uv} = 1$$
- path preservation:  $\forall_{u \neq s, t_p, t_c} \sum_v y_{vu} = \sum_v y_{uv}$
- path  $\neq$  cycle bounds:  $\sum_v y_{vt_{\{p,c\}}} \leq \sigma_{\{p,c\}}$
- forbid cycles (row generation via callback):
 
$$\forall \text{ cycle } C: \sum_{uv \in C} (y_{uv} - y_{ut_c}) < |C|$$
- Objective:
 
$$\max \sum_{e \in E} x_{\{u,v\}} \cdot \omega(e)$$
- cycle consistency:
 
$$\forall_u y_{ut_c} \leq y_{su}$$

# Extension: Contig Jumps

Mean read length: 70Bp

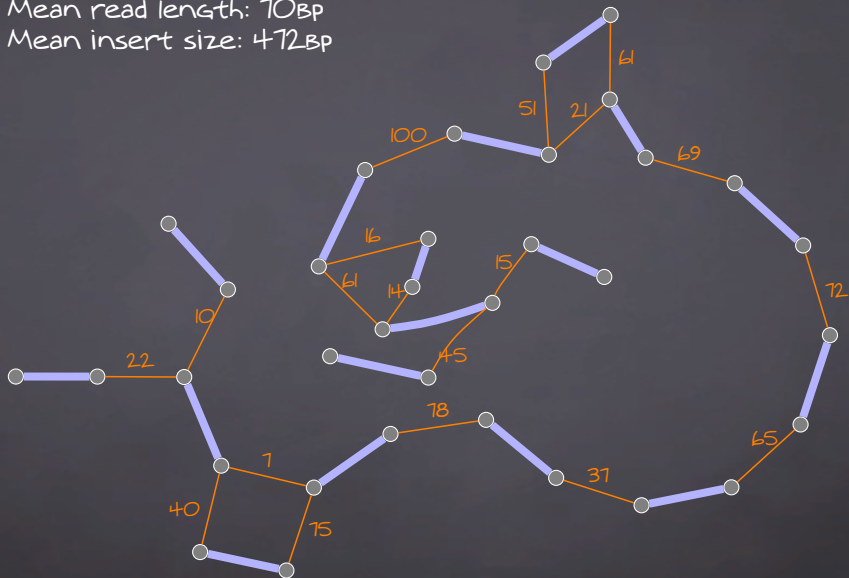
Mean insert size: 472Bp



# Extension: Contig Jumps

Mean read length: 70Bp

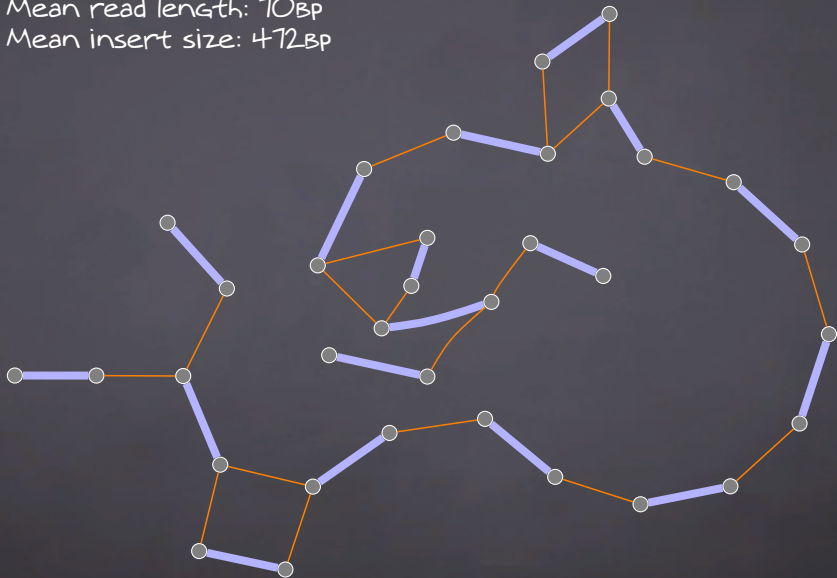
Mean insert size: 472Bp



# Extension: Contig Jumps

Mean read length: 70Bp

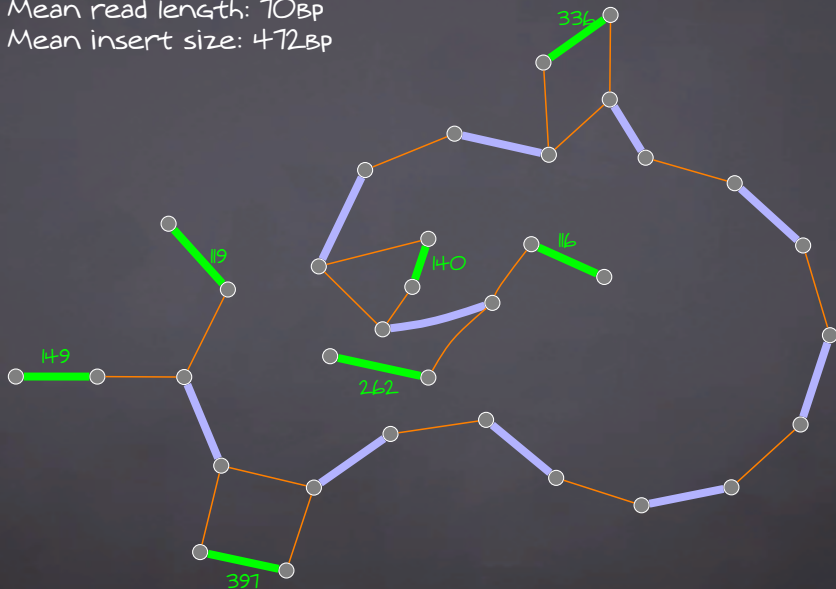
Mean insert size: 472Bp



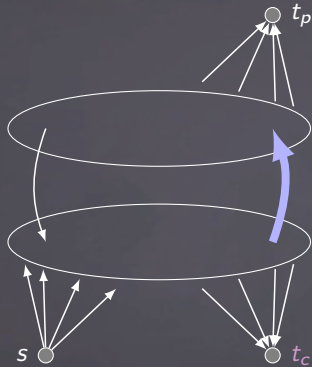
# Extension: Contig Jumps

Mean read length: 70Bp

Mean insert size: 472Bp

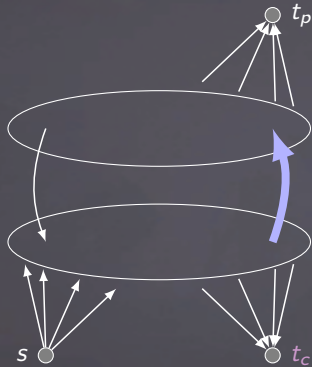


# Integer Linear Program Formulation

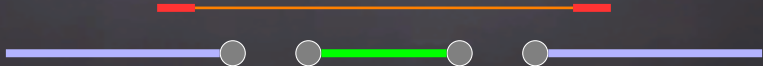


- chromosomes = disjoint  $s$ - $\{t_p, t_c\}$ -paths
- Bin. variables  $y_{uv} = 1 \Leftrightarrow u \rightarrow v$  used  
 $x_{\{u,v\}} = y_{uv} + y_{vu}$
- force contigs:  $\forall_{uv \in M} x_{uv} = 1$
- path preservation:  $\forall_{u \neq s, t_p, t_c} \sum_v y_{vu} = \sum_v y_{uv}$
- path  $\neq$  cycle bounds:  $\sum_v y_{vt_{\{p,c\}}} \leq \sigma_{\{p,c\}}$
- forbid cycles (row generation via callback):  
 $\forall \text{ cycle } C: \sum_{uv \in C} (y_{uv} - y_{ut_c}) < |C|$
- Objective:  $\max \sum_{e \in E} x_{\{u,v\}} \cdot \omega(e)$
- cycle consistency:  $\forall_u y_{ut_c} \leq y_{su}$

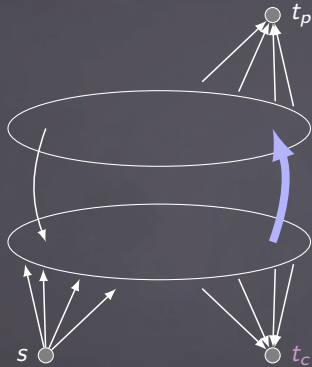
# Integer Linear Program Formulation



- chromosomes = disjoint  $s$ - $\{t_p, t_c\}$ -paths
- Bin. variables  $y_{uv} = 1 \Leftrightarrow u \rightarrow v$  used  
 $x_{\{u,v\}} = y_{uv} + y_{vu}$
- force contigs:  $\forall_{uv \in M} x_{uv} = 1$
- path preservation:  $\forall_{u \neq s, t_p, t_c} \sum_v y_{vu} = \sum_v y_{uv}$
- path  $\neq$  cycle bounds:  $\sum_v y_{vt_{\{p,c\}}} \leq \sigma_{\{p,c\}}$
- forbid cycles (row generation via callback):  
 $\forall \text{ cycle } C: \sum_{uv \in C} (y_{uv} - y_{ut_c}) < |C|$
- Objective:  $\max \sum_{e \in E} x_{\{u,v\}} \cdot \omega(e)$
- cycle consistency:  $\forall_u y_{ut_c} \leq y_{su}$



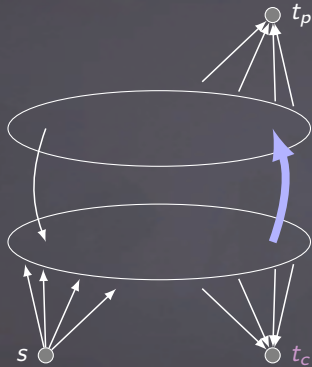
# Integer Linear Program Formulation



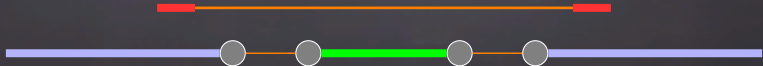
- chromosomes = disjoint  $s$ - $\{t_p, t_c\}$ -paths
- Bin. variables  $y_{uv} = 1 \Leftrightarrow u \rightarrow v$  used  
 $x_{\{u,v\}} = y_{uv} + y_{vu}$   
 $\forall_{uv \in M} x_{uv} = 1$
- force contigs:  $\forall_{u \neq s, t_p, t_c} \sum_v y_{vu} = \sum_v y_{uv}$
- path  $\neq$  cycle bounds:  $\sum_v y_{vt_{\{p,c\}}} \leq \sigma_{\{p,c\}}$
- forbid cycles (row generation via callback):  
 $\forall \text{ cycle } C: \sum_{uv \in C} (y_{uv} - y_{ut_c}) < |C|$
- Objective:  $\max \sum_{e \in E} x_{\{u,v\}} \cdot \omega(e)$
- cycle consistency:  $\forall_u y_{ut_c} \leq y_{su}$



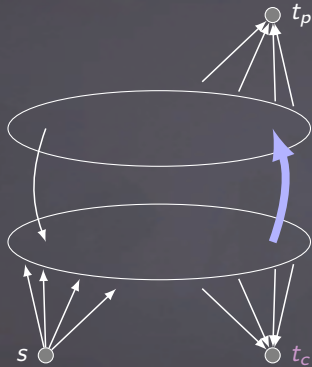
# Integer Linear Program Formulation



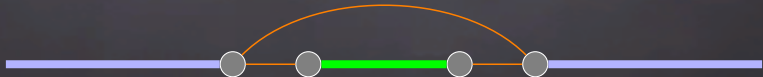
- chromosomes = disjoint  $s$ - $\{t_p, t_c\}$ -paths
- Bin. variables  $y_{uv} = 1 \Leftrightarrow u \rightarrow v$  used  
 $x_{\{u,v\}} = y_{uv} + y_{vu}$
- force contigs:  $\forall_{uv \in M} x_{uv} = 1$
- path preservation:  $\forall_{u \neq s, t_p, t_c} \sum_v y_{vu} = \sum_v y_{uv}$
- path  $\neq$  cycle bounds:  $\sum_v y_{vt_{\{p,c\}}} \leq \sigma_{\{p,c\}}$
- forbid cycles (row generation via callback):  
 $\forall \text{ cycle } C: \sum_{uv \in C} (y_{uv} - y_{ut_c}) < |C|$
- Objective:  $\max \sum_{e \in E} x_{\{u,v\}} \cdot \omega(e)$
- cycle consistency:  $\forall_u y_{ut_c} \leq y_{su}$



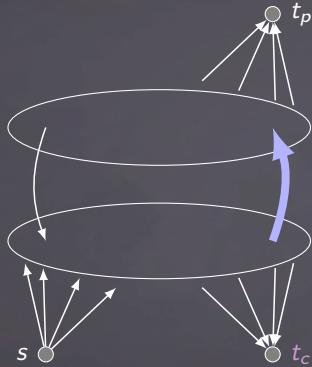
# Integer Linear Program Formulation



- chromosomes = disjoint  $s$ - $\{t_p, t_c\}$ -paths
- Bin. variables  $y_{uv} = 1 \Leftrightarrow u \rightarrow v$  used  
 $x_{\{u,v\}} = y_{uv} + y_{vu}$
- force contigs:  $\forall_{uv \in M} x_{uv} = 1$
- path preservation:  $\forall_{u \neq s, t_p, t_c} \sum_v y_{vu} = \sum_v y_{uv}$
- path  $\neq$  cycle bounds:  $\sum_v y_{vt_{\{p,c\}}} \leq \sigma_{\{p,c\}}$
- forbid cycles (row generation via callback):  
 $\forall \text{ cycle } C: \sum_{uv \in C} (y_{uv} - y_{ut_c}) < |C|$
- Objective:  $\max \sum_{e \in E} x_{\{u,v\}} \cdot \omega(e)$
- cycle consistency:  $\forall_u y_{ut_c} \leq y_{su}$



# Integer Linear Program Formulation



- chromosomes = disjoint  $s$ - $\{t_p, t_c\}$ -paths
- Bin. variables  $y_{uv} = 1 \Leftrightarrow u \rightarrow v$  used  
 $x_{\{u,v\}} = y_{uv} + y_{vu}$
- force contigs:  $\forall_{uv \in M} x_{uv} = 1$
- path preservation:  $\forall_{u \neq s, t_p, t_c} \sum_v y_{vu} = \sum_v y_{uv}$
- path  $\neq$  cycle bounds:  $\sum_v y_{vt_{\{p,c\}}} \leq \sigma_{\{p,c\}}$
- forbid cycles (row generation via callback):  
 $\forall$  cycle  $C$ :  $\sum_{uv \in C} (y_{uv} - y_{ut_c}) < |C|$
- Objective:  $\max \sum_{e \in E} x_{\{u,v\}} \cdot \omega(e)$
- cycle consistency:  $\forall_u y_{ut_c} \leq y_{su}$

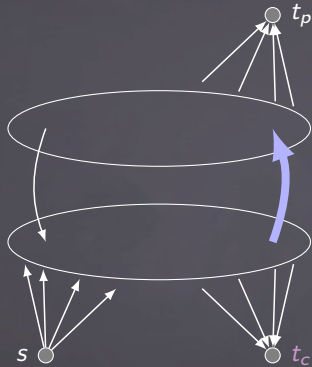
## Jump Mechanics

for each non-contig  $uv$ ,

1. introduce a variable  $z_{uv}$
2. construct "jump network" between  $u$  and  $v$  that fits in the gap
3. add  $z_{uv}$  to  $x_{\{u,v\}}$

extra: preprocess instance to finish incomplete jumps

# Integer Linear Program Formulation



- chromosomes = disjoint  $s-\{t_p, t_c\}$ -paths
- Bin. variables  
 $y_{uv} = 1 \Leftrightarrow u \rightarrow v$  used  
 $x_{\{u,v\}} = y_{uv} + y_{vu} + z_{uv} + z_{vu}$
- force contigs:  $\forall_{uv \in M} x_{uv} = 1$
- path preservation:  $\forall_{u \neq s, t_p, t_c} \sum_v y_{vu} = \sum_v y_{uv}$
- path  $\neq$  cycle bounds:  $\sum_v y_{vt_{\{p,c\}}} \leq \sigma_{\{p,c\}}$
- forbid cycles (row generation via callback):  
 $\forall \text{ cycle } C: \sum_{uv \in C} (y_{uv} - y_{ut_c}) < |C|$
- Objective:  $\max \sum_{e \in E} x_{\{u,v\}} \cdot \omega(e)$
- cycle consistency:  $\forall_u y_{ut_c} \leq y_{su}$
- jump mechanics

## Jump Mechanics

for each non-contig  $uv$ ,

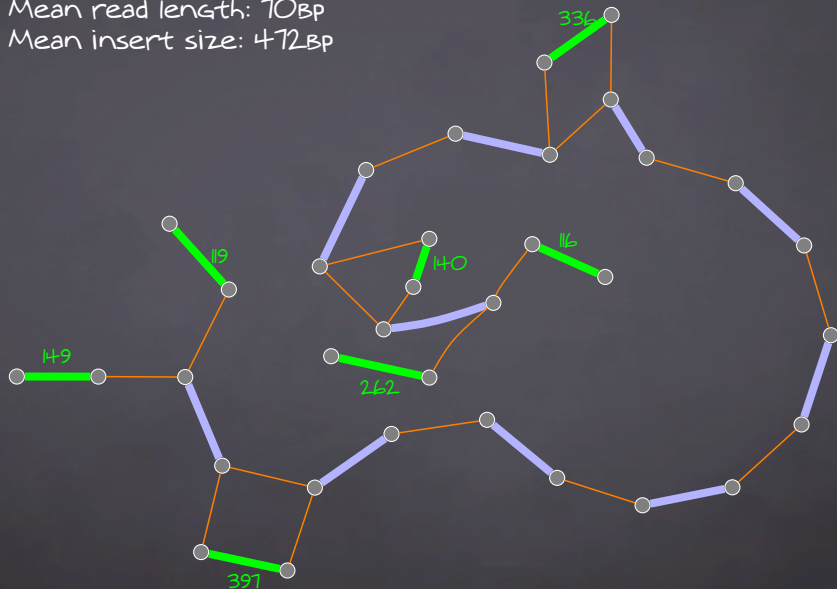
1. introduce a variable  $z_{uv}$
2. construct "jump network" between  $u$  and  $v$  that fits in the gap
3. add  $z_{uv}$  to  $x_{\{u,v\}}$

extra: preprocess instance to finish incomplete jumps

# Extension: Contig Jumps

Mean read length: 70Bp

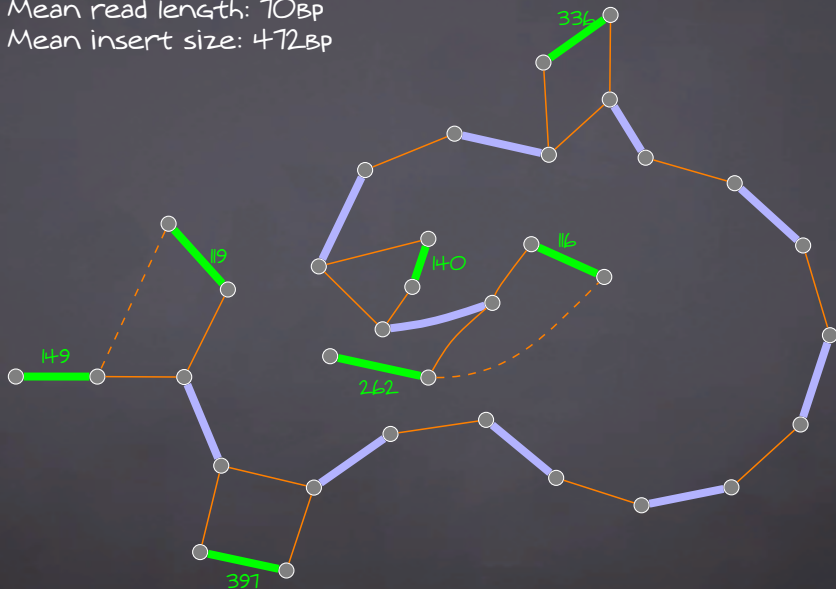
Mean insert size: 472Bp



# Extension: Contig Jumps

Mean read length: 70Bp

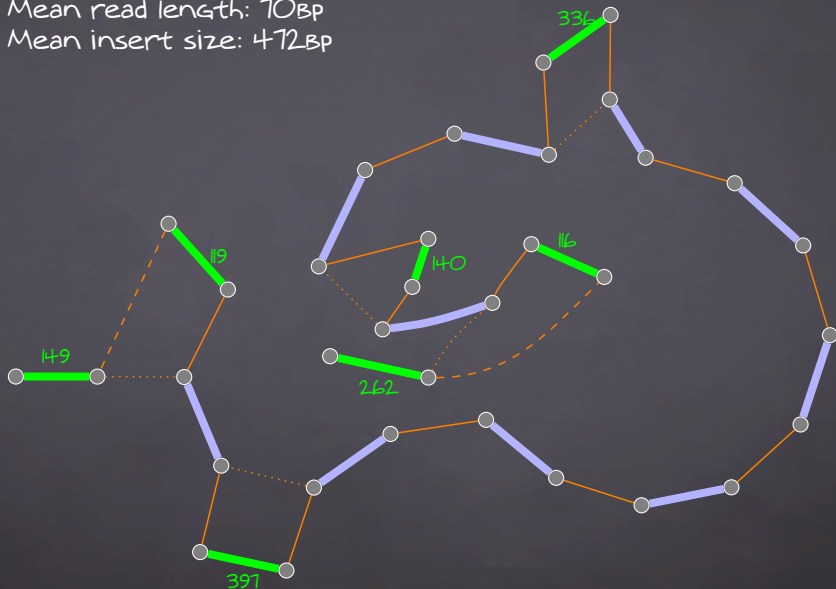
Mean insert size: 472Bp



# Extension: Contig Jumps

Mean read length: 70Bp

Mean insert size: 472Bp



## ILP Extension: Multiplicities

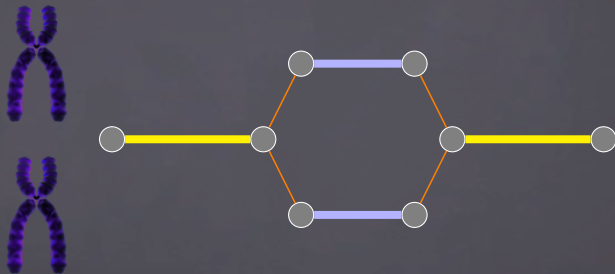


GGTGCGAGAGAGGTCATGGATTGCAACGA

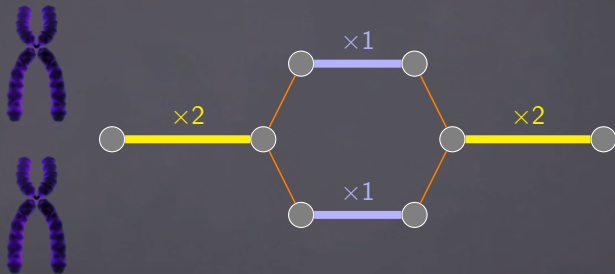


GGTGCGAGAGGCCACTCCAATTGCAACGA

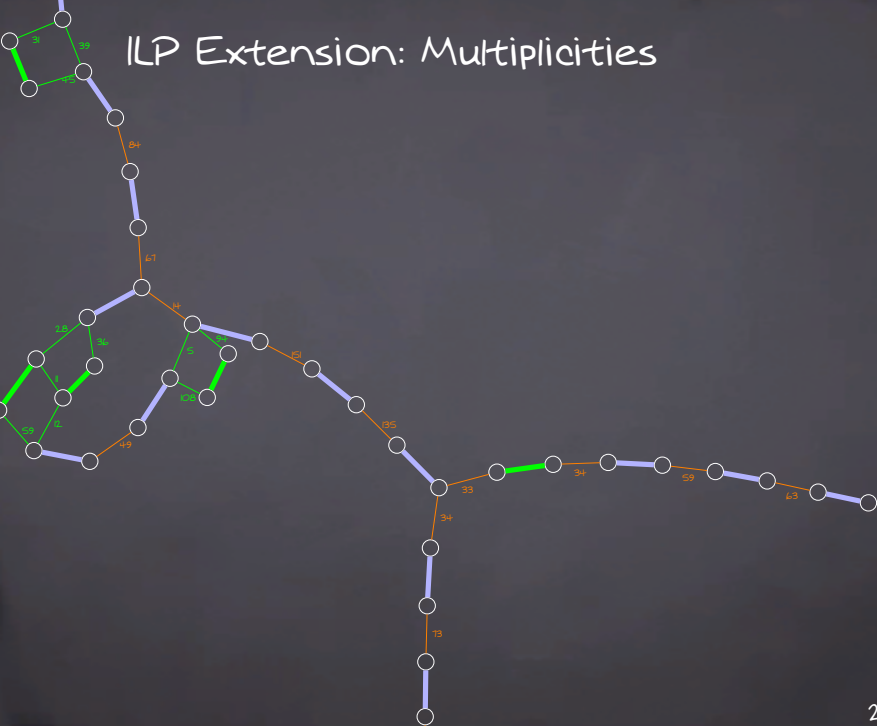
# ILP Extension: Multiplicities



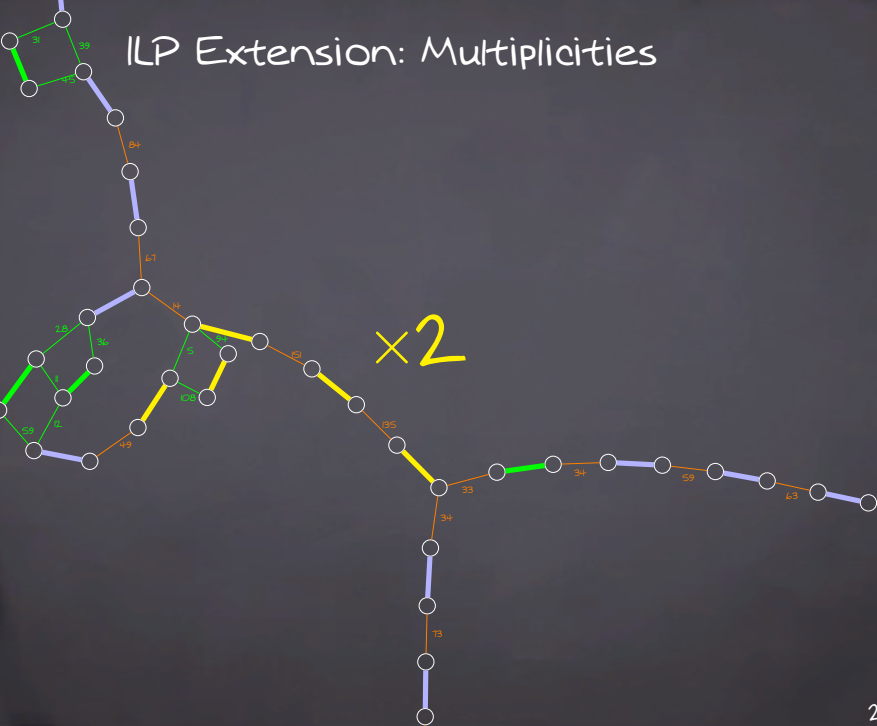
# ILP Extension: Multiplicities



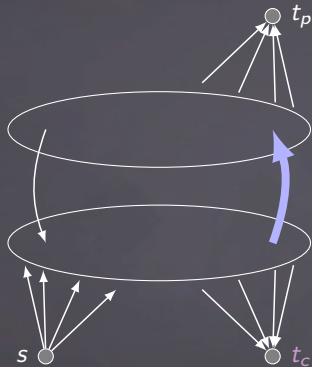
# ILP Extension: Multiplicities



# ILP Extension: Multiplicities



# Integer Linear Program Formulation

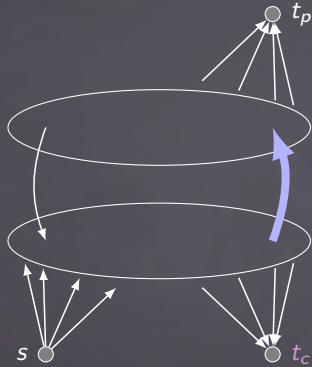


- chromosomes = disjoint  $s$ - $\{t_p, t_c\}$ -paths
- Bin. variables  $y_{uv} = 1 \Leftrightarrow u \rightarrow v$  used  
 $x_{\{u,v\}} = y_{uv} + y_{vu} + z_{uv} + z_{vu}$
- force contigs:  $\forall_{uv \in M} x_{uv} = 1$
- path preservation:  $\forall_{u \neq s, t_p, t_c} \sum_v y_{vu} = \sum_v y_{uv}$
- path  $\neq$  cycle bounds:  $\sum_v y_{vt_{\{p,c\}}} \leq \sigma_{\{p,c\}}$
- forbid cycles (row generation via callback):  
 $\forall \text{ cycle } C: \sum_{uv \in C} (y_{uv} - y_{ut_c}) < |C|$
- Objective:  $\max \sum_{e \in E} x_{\{u,v\}} \cdot \omega(e)$
- cycle consistency:  $\forall_u y_{ut_c} \leq y_{su}$
- jump mechanics

## Multiplicities

1. make  $y_{uv}, x_{\{u,v\}}$  integers in domain  $[0, m(\{u, v\})]$
2. change callback

# Integer Linear Program Formulation



- chromosomes = disjoint  $s-\{t_p, t_c\}$ -paths
- **int. variables**  $y_{uv} = \ell \Leftrightarrow u \rightarrow v$  used  $\ell$  times  
 $x_{\{u,v\}} = y_{uv} + y_{vu} + z_{uv} + z_{vu}$
- force contigs:  $\forall_{uv \in M} x_{uv} \geq 1$
- path preservation:  $\forall_{u \neq s, t_p, t_c} \sum_v y_{vu} = \sum_v y_{uv}$
- path  $\neq$  cycle bounds:  $\sum_v y_{vt_{\{p,c\}}} \leq \sigma_{\{p,c\}}$
- forbid cycles (row generation via callback):  
 $\forall \text{ cycle } C: \sum_{uv \in C} y_{uv} \leq |C| \cdot m_{\max} \cdot \sum_{u \in C, v \notin C} y_{uv}$
- **Objective:**  $\max \sum_{e \in E} x_{\{u,v\}} \cdot \omega(e)$
- cycle consistency:  $\forall_u y_{ut_c} \leq y_{su}$
- **jump mechanics !!!**

## Multiplicities

1. make  $y_{uv}, x_{\{u,v\}}$  integers in domain  $[0, m(\{u, v\})]$
2. **change callback**

# Linearization of Solutions

## Problem

no unique chromosome-configuration explaining solution



# Linearization of Solutions

## Problem

no unique chromosome-configuration explaining solution



# Linearization of Solutions

## Problem

no unique chromosome-configuration explaining solution



# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

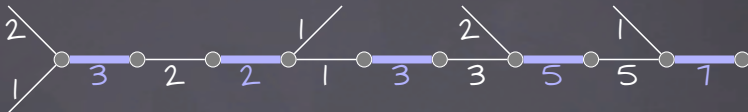
(=alt. path of uniform multiplicity  $\mu \nmid$  each end incident to non-contig  $< \mu$ )

# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(= alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )



# Linearization of Solutions

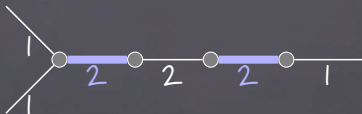
## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(= alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proof

" $\Rightarrow$ ": contraposition; let  $p$  = ambiguous path



# Linearization of Solutions

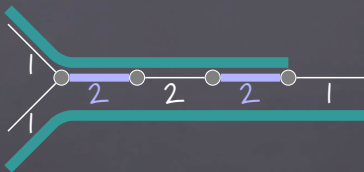
## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(= alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proof

" $\Rightarrow$ ": contraposition; let  $p$  = ambiguous path



# Linearization of Solutions

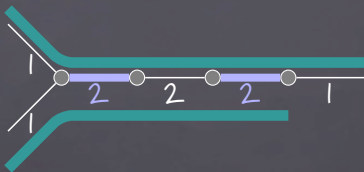
## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(= alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proof

" $\Rightarrow$ ": contraposition; let  $p$  = ambiguous path



# Linearization of Solutions

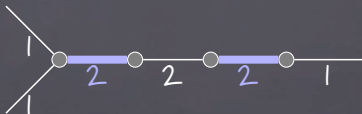
## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(= alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proof

" $\Rightarrow$ ": contraposition; let  $p$  = ambiguous path



$\leadsto (G, \mathcal{M}, m)$  not uniquely linearizable

# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(=alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proof

" $\Leftarrow$ ": let  $(G, \mathcal{M}, m)$  be free of ambiguous paths

Reduction (does not decrease number of linearizations):

# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(= alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proof

" $\Leftarrow$ ": let  $(G, \mathcal{M}, m)$  be free of ambiguous paths

Reduction (does not decrease number of linearizations):



# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(=alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proof

" $\Leftarrow$ ": let  $(G, \mathcal{M}, m)$  be free of ambiguous paths

Reduction (does not decrease number of linearizations):



# Linearization of Solutions

## Theorem

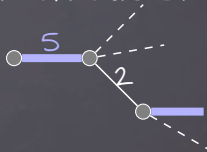
$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(= alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proof

" $\Leftarrow$ ": let  $(G, \mathcal{M}, m)$  be free of ambiguous paths

Reduction (does not decrease number of linearizations):



# Linearization of Solutions

## Theorem

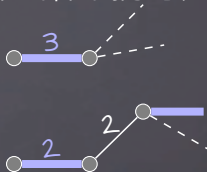
$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(= alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proof

" $\Leftarrow$ ": let  $(G, \mathcal{M}, m)$  be free of ambiguous paths

Reduction (does not decrease number of linearizations):



# Linearization of Solutions

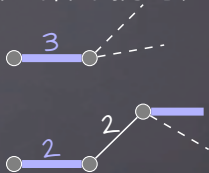
## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"  
(= alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proof

" $\Leftarrow$ ": let  $(G, \mathcal{M}, m)$  be free of ambiguous paths

Reduction (does not decrease number of linearizations):



$\leadsto$  result is collection of alternating paths  $\neq$  cycles

# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(=alt. path of uniform multiplicity  $\mu \nmid$  each end incident to non-contig  $< \mu$ )

## Proposals

# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(=alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proposals

- I. decide arbitrarily

# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(=alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proposals

1. decide arbitrarily  $\rightsquigarrow$  missassembly

# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(=alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proposals

1. decide arbitrarily  $\rightsquigarrow$  missassembly
2. isolate each ambiguity

# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(=alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proposals

1. decide arbitrarily  $\rightsquigarrow$  missassembly
2. isolate each ambiguity  $\rightsquigarrow$  information loss

# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(=alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proposals

1. decide arbitrarily  $\rightsquigarrow$  missassembly
2. isolate each ambiguity  $\rightsquigarrow$  information loss
3. cut as few ends as possible

# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(=alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proposals

1. decide arbitrarily  $\rightsquigarrow$  missassembly
2. isolate each ambiguity  $\rightsquigarrow$  information loss
3. cut as few ends as possible  $\rightsquigarrow$  computationally hard

# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(=alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proposals

1. decide arbitrarily  $\rightsquigarrow$  missassembly
2. isolate each ambiguity  $\rightsquigarrow$  information loss
3. cut as few ends as possible  $\rightsquigarrow$  computationally hard
4. cut as few multiplicities as possible

# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(=alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proposals

1. decide arbitrarily  $\rightsquigarrow$  missassembly
2. isolate each ambiguity  $\rightsquigarrow$  information loss
3. cut as few ends as possible  $\rightsquigarrow$  computationally hard
4. cut as few multiplicities as possible  $\rightsquigarrow$  computationally hard

# Linearization of Solutions

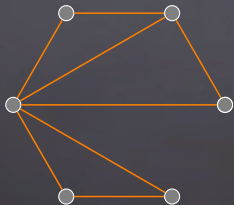
## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(= alt. path of uniform multiplicity  $\mu \nmid$  each end incident to non-contig  $< \mu$ )

## Proposals

1. decide arbitrarily  $\rightsquigarrow$  missassembly
2. isolate each ambiguity  $\rightsquigarrow$  information loss
3. cut as few ends as possible  $\rightsquigarrow$  computationally hard
4. cut as few multiplicities as possible  $\rightsquigarrow$  computationally hard



## Multiplicities

one  $\nmid$

#non-matching adj. to contig

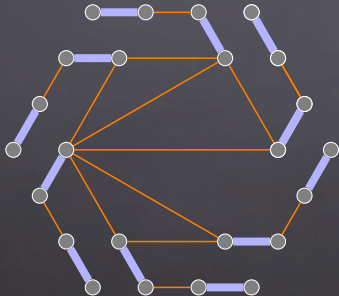
# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"  
(=alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proposals

1. decide arbitrarily  $\rightsquigarrow$  missassembly
2. isolate each ambiguity  $\rightsquigarrow$  information loss
3. cut as few ends as possible  $\rightsquigarrow$  computationally hard
4. cut as few multiplicities as possible  $\rightsquigarrow$  computationally hard



Multiplicities

one  $\neq$

#non-matching adj. to contig

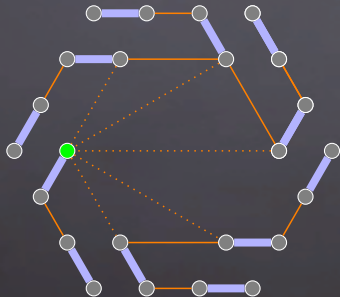
# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"  
(=alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proposals

1. decide arbitrarily  $\rightsquigarrow$  missassembly
2. isolate each ambiguity  $\rightsquigarrow$  information loss
3. cut as few ends as possible  $\rightsquigarrow$  computationally hard
4. cut as few multiplicities as possible  $\rightsquigarrow$  computationally hard



Multiplicities

one  $\neq$

#non-matching adj. to contig

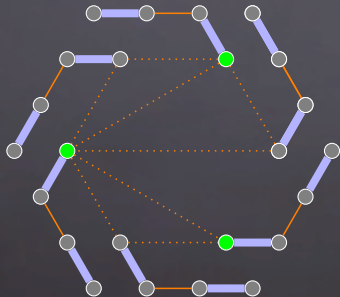
# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"  
(=alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proposals

1. decide arbitrarily  $\rightsquigarrow$  missassembly
2. isolate each ambiguity  $\rightsquigarrow$  information loss
3. cut as few ends as possible  $\rightsquigarrow$  computationally hard
4. cut as few multiplicities as possible  $\rightsquigarrow$  computationally hard



Multiplicities

one  $\neq$

#non-matching adj. to contig

# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(= alt. path of uniform multiplicity  $\mu \nmid$  each end incident to non-contig  $< \mu$ )

## Proposals

1. decide arbitrarily  $\rightsquigarrow$  missassembly
2. isolate each ambiguity  $\rightsquigarrow$  information loss
3. cut as few ends as possible  $\rightsquigarrow$  computationally hard
4. cut as few multiplicities as possible  $\rightsquigarrow$  computationally hard



# Linearization of Solutions

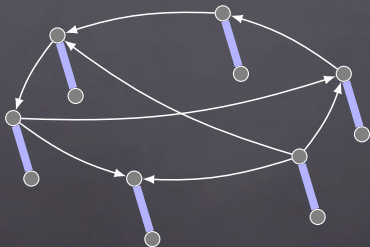
## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(= alt. path of uniform multiplicity  $\mu \nmid$  each end incident to non-contig  $< \mu$ )

## Proposals

1. decide arbitrarily  $\rightsquigarrow$  missassembly
2. isolate each ambiguity  $\rightsquigarrow$  information loss
3. cut as few ends as possible  $\rightsquigarrow$  computationally hard
4. cut as few multiplicities as possible  $\rightsquigarrow$  computationally hard



# Linearization of Solutions

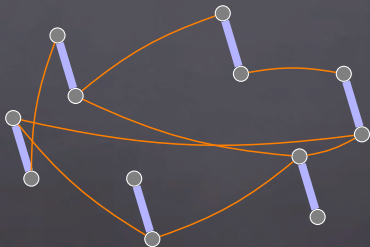
## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(= alt. path of uniform multiplicity  $\mu \neq$  each end incident to non-contig  $< \mu$ )

## Proposals

1. decide arbitrarily  $\rightsquigarrow$  missassembly
2. isolate each ambiguity  $\rightsquigarrow$  information loss
3. cut as few ends as possible  $\rightsquigarrow$  computationally hard
4. cut as few multiplicities as possible  $\rightsquigarrow$  computationally hard



# Linearization of Solutions

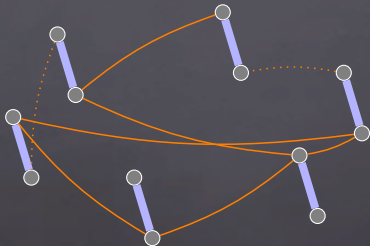
## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"

(= alt. path of uniform multiplicity  $\mu \nmid$  each end incident to non-contig  $< \mu$ )

## Proposals

1. decide arbitrarily  $\rightsquigarrow$  missassembly
2. isolate each ambiguity  $\rightsquigarrow$  information loss
3. cut as few ends as possible  $\rightsquigarrow$  computationally hard
4. cut as few multiplicities as possible  $\rightsquigarrow$  computationally hard



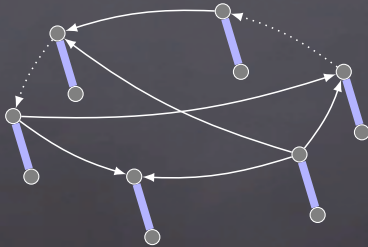
# Linearization of Solutions

## Theorem

$(G, \mathcal{M}, m)$  uniquely linearizable  $\Leftrightarrow$  no "ambiguous paths"  
(= alt. path of uniform multiplicity  $\mu \nmid$  each end incident to non-contig  $< \mu$ )

## Proposals

1. decide arbitrarily  $\rightsquigarrow$  missassembly
2. isolate each ambiguity  $\rightsquigarrow$  information loss
3. cut as few ends as possible  $\rightsquigarrow$  computationally hard
4. cut as few multiplicities as possible  $\rightsquigarrow$  computationally hard



# Conclusion

## What we saw

- 3-step sequencing technique:
  1. produce paired-end reads
  2. assemble reads to contigs
  3. scaffold contigs to chromosomes using read-pairings

# Conclusion

## What we saw

- 3-step sequencing technique:
  1. produce paired-end reads
  2. assemble reads to contigs
  3. scaffold contigs to chromosomes using read-pairings
- computationally hard problem for dense graphs with weights  $O(1)$
- no constant-factor approx or subexponential-time algorithm for linear quasi trees with weights  $O(1)$

# Conclusion

## What we saw

- 3-step sequencing technique:
  1. produce paired-end reads
  2. assemble reads to contigs
  3. scaffold contigs to chromosomes using read-pairings
- computationally hard problem for dense graphs with weights  $O(1)$
- no constant-factor approx or subexponential-time algorithm for linear quasi trees with weights  $O(1)$
- $O(n^2)$  time on unweighted cliques/co-bipartite/split

# Conclusion

## What we saw

- 3-step sequencing technique:
  1. produce paired-end reads
  2. assemble reads to contigs
  3. scaffold contigs to chromosomes using read-pairings
- computationally hard problem for dense graphs with weights  $O(1)$
- no constant-factor approx or subexponential-time algorithm for linear quasi trees with weights  $O(1)$
- $O(n^2)$  time on unweighted cliques/co-bipartite/split
- $O(n \cdot \sigma_p \cdot \sigma_c)$  time for constant treewidth

# Conclusion

## What we saw

- 3-step sequencing technique:
  1. produce paired-end reads
  2. assemble reads to contigs
  3. scaffold contigs to chromosomes using read-pairings
- computationally hard problem for dense graphs with weights  $O(1)$
- no constant-factor approx or subexponential-time algorithm for linear quasi trees with weights  $O(1)$
- $O(n^2)$  time on unweighted cliques/co-bipartite/split
- $O(n \cdot \sigma_p \cdot \sigma_c)$  time for constant treewidth
- 2-approximable in cliques/complete bipartite in  $O(n^3)$  time

# Conclusion

## What we saw

- 3-step sequencing technique:
  1. produce paired-end reads
  2. assemble reads to contigs
  3. scaffold contigs to chromosomes using read-pairings
- computationally hard problem for dense graphs with weights  $O(1)$
- no constant-factor approx or subexponential-time algorithm for linear quasi trees with weights  $O(1)$
- $O(n^2)$  time on unweighted cliques/co-bipartite/split
- $O(n \cdot \sigma_p \cdot \sigma_c)$  time for constant treewidth
- 2-approximable in cliques/complete bipartite in  $O(n^3)$  time
- $O(\sqrt{2}^n \text{poly}(n))$  time exact algorithm

# Conclusion

## What we saw

- 3-step sequencing technique:
  1. produce paired-end reads
  2. assemble reads to contigs
  3. scaffold contigs to chromosomes using read-pairings
- computationally hard problem for dense graphs with weights  $O(1)$
- no constant-factor approx or subexponential-time algorithm for linear quasi trees with weights  $O(1)$
- $O(n^2)$  time on unweighted cliques/co-bipartite/split
- $O(n \cdot \sigma_p \cdot \sigma_c)$  time for constant treewidth
- 2-approximable in cliques/complete bipartite in  $O(n^3)$  time
- $O(\sqrt{2}^n \text{poly}(n))$  time exact algorithm
- ILP formulation with **contig jumps**  $\neq$  **multiplicities**

# Conclusion

## What we saw

- 3-step sequencing technique:
  1. produce paired-end reads
  2. assemble reads to contigs
  3. scaffold contigs to chromosomes using read-pairings
- computationally hard problem for dense graphs with weights  $O(1)$
- no constant-factor approx or subexponential-time algorithm for linear quasi trees with weights  $O(1)$
- $O(n^2)$  time on unweighted cliques/co-bipartite/split
- $O(n \cdot \sigma_p \cdot \sigma_c)$  time for constant treewidth
- 2-approximable in cliques/complete bipartite in  $O(n^3)$  time
- $O(\sqrt{2}^n \text{poly}(n))$  time exact algorithm
- ILP formulation with **contig jumps**  $\neq$  **multiplicities**
- Linearization problem raised by multiplicities in solution

# Conclusion

## Outlook

- 3<sup>rd</sup> generation sequencing: PacBio, Oxford Nanopore produces long reads (10-15kbp), but error-prone  
→ correction using small reads?

# Conclusion

## Outlook

- 3<sup>rd</sup> generation sequencing: PacBio, Oxford Nanopore produces long reads (10-15kbp), but error-prone  
    ~> correction using small reads?
- Generally: multi-library scaffolding

# Conclusion

## Outlook

- 3<sup>rd</sup> generation sequencing: PacBio, Oxford Nanopore produces long reads (10-15kbp), but error-prone  
    ⇒ correction using small reads?
- Generally: multi-library scaffolding
- other sources for contig-connections (phylogenetic information?)

# Conclusion

## Outlook

- 3<sup>rd</sup> generation sequencing: PacBio, Oxford Nanopore produces long reads (10-15kbp), but error-prone  
    ↪ correction using small reads?
- Generally: multi-library scaffolding
- other sources for contig-connections (phylogenetic information?)
- Better parameters for Scaffolding and Scaffold Linearization  
    ↪ analyze practical instances

# Conclusion

## Outlook

- 3<sup>rd</sup> generation sequencing: PacBio, Oxford Nanopore produces long reads (10-15kbp), but error-prone  
    ~> correction using small reads?
- Generally: multi-library scaffolding
- other sources for contig-connections (phylogenetic information?)
- Better parameters for Scaffolding and Scaffold Linearization  
    ~> analyze practical instances
- approximation/heuristics for Scaffold Linearization

