# Advanced Dynamic Programming in CL:

# Theory, Algorithms, and Applications

$(S, 0, n)$

$w_0 \quad w_1 \quad \ldots \quad w_{n-1}$

Liang Huang
University of Pennsylvania

# A Little Bit of History...

# A Little Bit of History...

- Who invented Dynamic Programming?
  and when was it invented?

# A Little Bit of History...

- Who invented Dynamic Programming? and when was it invented?

- R. Bellman (1940s-50s)

- A. Viterbi (1967)

- E. Dijkstra (1959)

- Hart, Nilsson, and Raphael (1968)
  - Dijkstra => A* Algorithm

- D. Knuth (1977)
  - Dijkstra on Grammar (Hypergraph)

Richard Bellman

Andrew Viterbi

# A Little Bit of History...

- Who invented Dynamic Programming? and when was it invented?

- R. Bellman (1940s-50s)

- A. Viterbi (1967)

- E. Dijkstra (1959)

- Hart, Nilsson, and Raphael (1968)

  - Dijkstra => A* Algorithm

- D. Knuth (1977)

  - Dijkstra on Grammar (Hypergraph)

A. Turing

Richard Bellman

Andrew Viterbi

# Dynamic Programming

- Dynamic Programming is everywhere in NLP

  - Viterbi Algorithm for Hidden Markov Models

  - CKY Algorithm for Parsing and Machine Translation

  - Forward-Backward and Inside-Outside Algorithms

- Also everywhere in AI/ML

  - Reinforcement Learning, Planning (POMDP)

  - AI Search: Uniform-cost, A*, etc.

- This tutorial: a unified theoretical view of DP
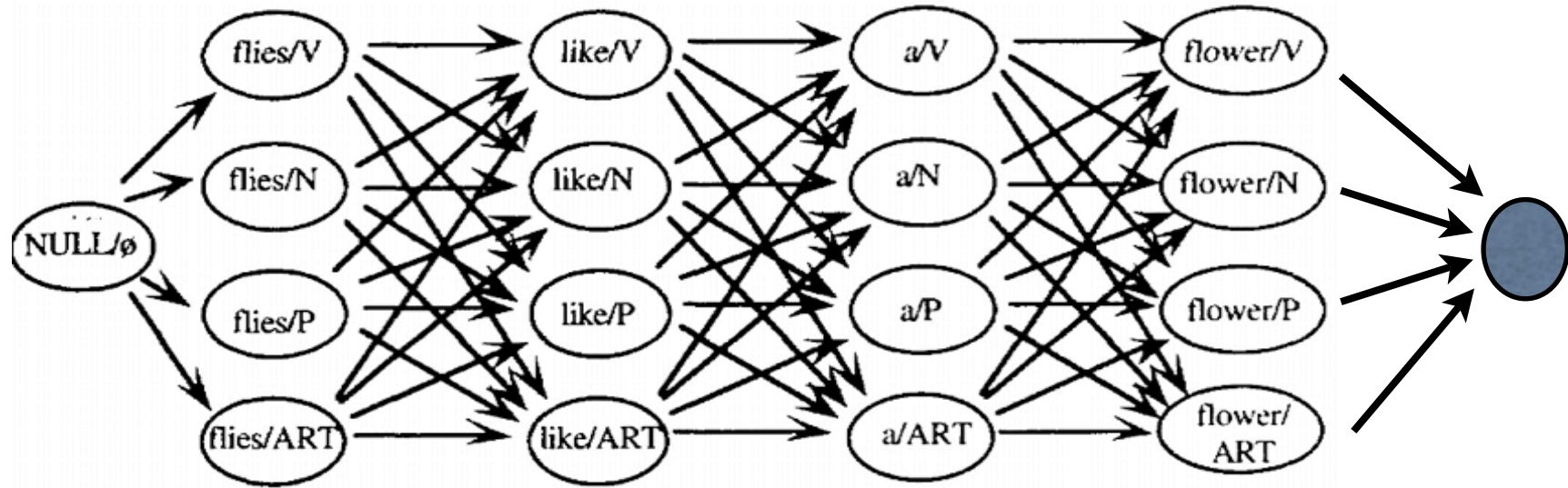
  - Focusing on *Optimization Problems*

# Review: DP Basics

- DP = Divide-and-Conquer + Two Principles:
  - [required] Optimal Subproblem Property
  - [recommended] Sharing of Common Subproblems

- Structure of the Search Space

  - Incremental
    - Graph
    - Knapsack, Edit Dist., Sequence Alignment

  - Branching
    - Hypergraph
    - Matrix-Chain, Polygon Triangulation, Optimal BST

# Two Dimensional Survey

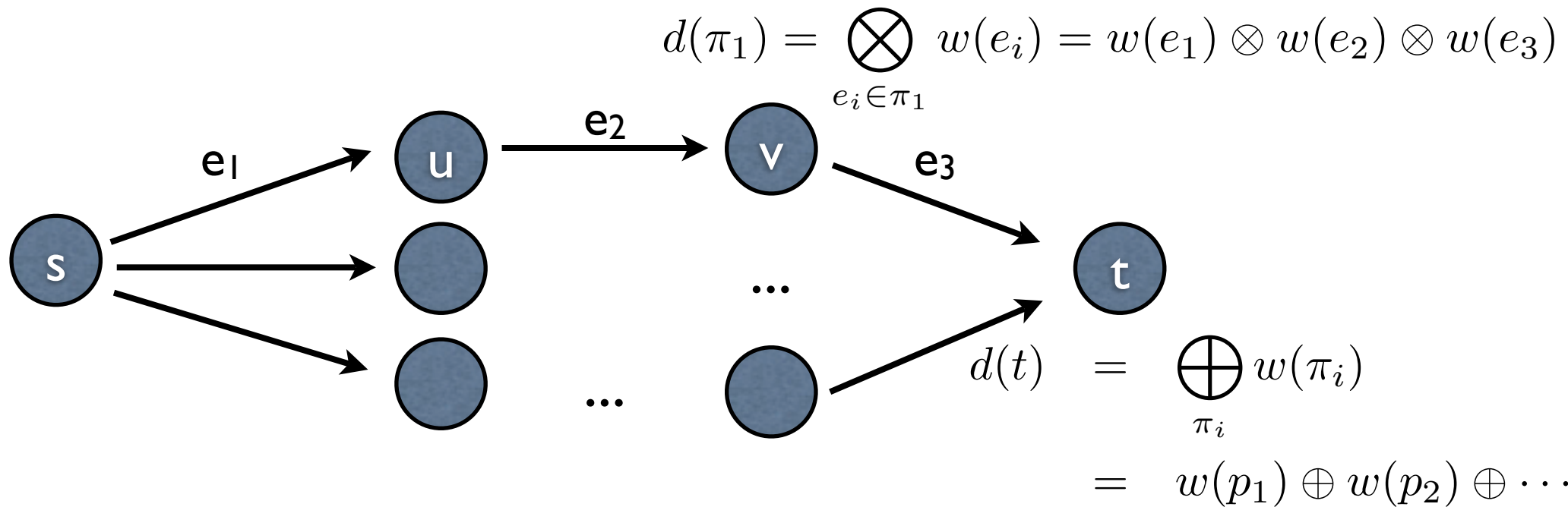| | topological (acyclic) | best-first (superior) |
|---|---|---|
| graphs with semirings | Viterbi | Dijkstra |
| hypergraphs with weight functions | Generalized Viterbi | Knuth |

search space

# Graphs in NLP

part-of-speech tagging



lattice in speech

# Semirings on Graphs

- in a weighted graph, we need two operators:

  - extension (multiplicative) and summary (additive)

  - the weight of a path is the product of edge weights

  - the weight of a vertex is the summary of path weights

$$d(\pi_1) = \bigotimes_{e_i \in \pi_1} w(e_i) = w(e_1) \otimes w(e_2) \otimes w(e_3)$$



$$d(t) = \bigoplus_{\pi_i} w(\pi_i)$$
$$= w(p_1) \oplus w(p_2) \oplus \cdots$$

# Semiring Definitions

A **monoid** is a triple $(A, \otimes, \overline{1})$ where

1. $\otimes$ is a closed <span style="color:red">associative binary operator</span> on the set $A$,

2. $\overline{1}$ is the <span style="color:red">identity element</span> for $\otimes$, i.e., for all $a \in A, a \otimes \overline{1} = \overline{1} \otimes a = a$.

A monoid is **commutative** if $\otimes$ is commutative.

# Semiring Definitions

A **monoid** is a triple $(A, \otimes, \overline{1})$ where

1. $\otimes$ is a closed associative binary operator on the set $A$,

2. $\overline{1}$ is the identity element for $\otimes$, i.e., for all $a \in A, a \otimes \overline{1} = \overline{1} \otimes a = a$.

A monoid is **commutative** if $\otimes$ is commutative.

$$([0, 1], +, 0)$$
$$([0, 1], \times, 1)$$
$$([0, 1], \max, 0)$$

# Semiring Definitions

A **monoid** is a triple $(A, \otimes, \overline{1})$ where

    1. $\otimes$ is a closed <span style="color:red">associative binary operator</span> on the set $A$,

    2. $\overline{1}$ is the <span style="color:red">identity element</span> for $\otimes$, i.e., for all $a \in A, a \otimes \overline{1} = \overline{1} \otimes a = a$.

A monoid is **commutative** if $\otimes$ is commutative.

$$([0, 1], +, 0)$$
$$([0, 1], \times, 1)$$
$$([0, 1], \max, 0)$$

A **semiring** is a 5-tuple $R = (A, \oplus, \otimes, \overline{0}, \overline{1})$ such that

1. $(A, \oplus, \overline{0})$ is a commutative monoid.

2. $(A, \otimes, \overline{1})$ is a monoid.

3. $\otimes$ distributes over $\oplus$: for all $a, b, c$ in $A$,

$$(a \oplus b) \otimes c = (a \otimes c) \oplus (b \otimes c),$$

$$c \otimes (a \oplus b) = (c \otimes a) \oplus (c \otimes b).$$

4. $\overline{0}$ is an **annihilator** for $\otimes$: for all $a$ in $A$, $\overline{0} \otimes a = a \otimes \overline{0} = \overline{0}$.

# Semiring Definitions

A **monoid** is a triple $(A, \otimes, \overline{1})$ where

1. $\otimes$ is a closed associative binary operator on the set $A$,

2. $\overline{1}$ is the identity element for $\otimes$, i.e., for all $a \in A$, $a \otimes \overline{1} = \overline{1} \otimes a = a$.

A monoid is **commutative** if $\otimes$ is commutative.

$$([0, 1], +, 0)$$
$$([0, 1], \times, 1)$$
$$([0, 1], \max, 0)$$

A **semiring** is a 5-tuple $R = (A, \oplus, \otimes, \overline{0}, \overline{1})$ such that

1. $(A, \oplus, \overline{0})$ is a commutative monoid.

2. $(A, \otimes, \overline{1})$ is a monoid.

$$([0, 1], \max, \times, 0, 1)$$
$$([0, 1], +, \times, 0, 1)$$

3. $\otimes$ distributes over $\oplus$: for all $a, b, c$ in $A$,

$$(a \oplus b) \otimes c = (a \otimes c) \oplus (b \otimes c),$$

$$c \otimes (a \oplus b) = (c \otimes a) \oplus (c \otimes b).$$

4. $\overline{0}$ is an **annihilator** for $\otimes$: for all $a$ in $A$, $\overline{0} \otimes a = a \otimes \overline{0} = \overline{0}$.

# Semiring Definitions

A **monoid** is a triple $(A, \otimes, \overline{1})$ where

1. $\otimes$ is a closed associative binary operator on the set $A$,

2. $\overline{1}$ is the identity element for $\otimes$, i.e., for all $a \in A$, $a \otimes \overline{1} = \overline{1} \otimes a = a$.

A monoid is **commutative** if $\otimes$ is commutative.

$$([0, 1], +, 0)$$
$$([0, 1], \times, 1)$$
$$([0, 1], \max, 0)$$

A **semiring** is a 5-tuple $R = (A, \oplus, \otimes, \overline{0}, \overline{1})$ such that

1. $(A, \oplus, \overline{0})$ is a commutative monoid.

2. $(A, \otimes, \overline{1})$ is a monoid.

3. $\otimes$ distributes over $\oplus$: for all $a, b, c$ in $A$,

$$([0, 1], \max, \times, 0, 1) \checkmark$$
$$([0, 1], +, \times, 0, 1) \; ✗$$

$$(a \oplus b) \otimes c = (a \otimes c) \oplus (b \otimes c),$$

$$c \otimes (a \oplus b) = (c \otimes a) \oplus (c \otimes b).$$

4. $\overline{0}$ is an **annihilator** for $\otimes$: for all $a$ in $A$, $\overline{0} \otimes a = a \otimes \overline{0} = \overline{0}$.

# Examples

| Semiring | Set | $\oplus$ | $\otimes$ | $\overline{0}$ | $\overline{1}$ | intuition/application |
|---|---|---|---|---|---|---|
| Boolean | $\{0, 1\}$ | $\vee$ | $\wedge$ | 0 | 1 | logical deduction, recognition |
| Viterbi | $[0, 1]$ | max | $\times$ | 0 | 1 | prob. of the best derivation |
| Inside | $\mathbb{R}^+ \cup \{+\infty\}$ | $+$ | $\times$ | 0 | 1 | prob. of a string |
| Real | $\mathbb{R} \cup \{+\infty\}$ | **min** | $+$ | $+\infty$ | 0 | shortest-distance |
| Tropical | $\mathbb{R}^+ \cup \{+\infty\}$ | **min** | $+$ | $+\infty$ | 0 | with non-negative weights |
| Counting | $\mathbb{N}$ | $+$ | $\times$ | 0 | 1 | number of paths |

# Ordering

# Ordering

- ## idempotent

A semiring $(A, \oplus, \otimes, \overline{0}, \overline{1})$ is **idempotent** if for all $a$ in $A$, $a \oplus a = a$.

# Ordering

- ## idempotent

A semiring $(A, \oplus, \otimes, \overline{0}, \overline{1})$ is **idempotent** if for all $a$ in $A$, $a \oplus a = a$.

- ## comparison

$(a \leq b) \Leftrightarrow (a \oplus b = a)$ defines a partial ordering.

- ## examples: boolean, viterbi, tropical, real, ...

$$(\{0, 1\}, \vee, \wedge, 0, 1) \qquad (\mathbb{R}^+ \cup \{+\infty\}, \mathbf{min}, +, +\infty, 0)$$

$$([0, 1], \max, \otimes, 0, 1) \qquad (\mathbb{R} \cup \{+\infty\}, \mathbf{min}, +, +\infty, 0)$$

# Ordering

- **idempotent**

  A semiring $(A, \oplus, \otimes, \overline{0}, \overline{1})$ is **idempotent** if for all $a$ in $A$, $a \oplus a = a$.

- **comparison**

  $(a \leq b) \Leftrightarrow (a \oplus b = a)$ defines a partial ordering.

- **examples: boolean, viterbi, tropical, real, ...**

$$(\{0, 1\}, \vee, \wedge, 0, 1) \qquad (\mathbb{R}^+ \cup \{+\infty\}, \mathbf{min}, +, +\infty, 0)$$

$$([0, 1], \max, \otimes, 0, 1) \qquad (\mathbb{R} \cup \{+\infty\}, \mathbf{min}, +, +\infty, 0)$$

- **total-order for optimization problems**

  A semiring is **totally-ordered** if $\oplus$ defines a total ordering.

- **examples: all of the above**

# Monotonicity

# Monotonicity

- monotonicity

# Monotonicity

- ## monotonicity

Let $K = (A, \oplus, \otimes, \overline{0}, \overline{1})$ be a semiring, and $\leq$ a partial ordering over $A$. We say $K$ is **monotonic** if for all $a, b, c \in A$

$$(a \leq b) \Rightarrow (a \otimes c \leq b \otimes c) \qquad (a \leq b) \Rightarrow (c \otimes a \leq c \otimes b)$$

# Monotonicity

- ## monotonicity

Let $K = (A, \oplus, \otimes, \overline{0}, \overline{1})$ be a semiring, and $\leq$ a partial ordering over $A$. We say $K$ is **monotonic** if for all $a, b, c \in A$

$$(a \leq b) \Rightarrow (a \otimes c \leq b \otimes c) \qquad (a \leq b) \Rightarrow (c \otimes a \leq c \otimes b)$$

- ## optimal substructure in dynamic programming

# Monotonicity

- **monotonicity**

Let $K = (A, \oplus, \otimes, \overline{0}, \overline{1})$ be a semiring, and $\leq$ a partial ordering over $A$. We say $K$ is **monotonic** if for all $a, b, c \in A$

$$(a \leq b) \Rightarrow (a \otimes c \leq b \otimes c) \qquad (a \leq b) \Rightarrow (c \otimes a \leq c \otimes b)$$

- **optimal substructure** in dynamic programming

A: $b \otimes c$

B: $b$

C: $c$

# Monotonicity

- **monotonicity**

Let $K = (A, \oplus, \otimes, \overline{0}, \overline{1})$ be a semiring, and $\leq$ a partial ordering over $A$. We say $K$ is **monotonic** if for all $a, b, c \in A$

$$(a \leq b) \Rightarrow (a \otimes c \leq b \otimes c) \qquad (a \leq b) \Rightarrow (c \otimes a \leq c \otimes b)$$

- **optimal substructure** in dynamic programming

A: $b \otimes c$

B: $b$    C: $c$

A: $b' \otimes c$
$\leq b \otimes c$

B: $b' \leq b$    C: $c$

# Monotonicity

- **monotonicity**

Let $K = (A, \oplus, \otimes, \overline{0}, \overline{1})$ be a semiring, and $\leq$ a partial ordering over $A$. We say $K$ is **monotonic** if for all $a, b, c \in A$

$$(a \leq b) \Rightarrow (a \otimes c \leq b \otimes c) \qquad (a \leq b) \Rightarrow (c \otimes a \leq c \otimes b)$$

- **optimal substructure** in dynamic programming



A: $b \otimes c$

B: $b$    C: $c$

A: $b' \otimes c$
$\leq b \otimes c$

B: $b' \leq b$    C: $c$

- idempotent => monotone (from distributivity)

# Monotonicity

- **monotonicity**

Let $K = (A, \oplus, \otimes, \overline{0}, \overline{1})$ be a semiring, and $\leq$ a partial ordering over $A$. We say $K$ is **monotonic** if for all $a, b, c \in A$

$$(a \leq b) \Rightarrow (a \otimes c \leq b \otimes c) \qquad (a \leq b) \Rightarrow (c \otimes a \leq c \otimes b)$$

- **optimal substructure** in dynamic programming



- idempotent => monotone (from distributivity)

  - (a+b)⊗c = (a⊗c)+(b⊗c); if a≤b, (a⊗c)=(a⊗c)+(b⊗c)

# Monotonicity

- **monotonicity**

Let $K = (A, \oplus, \otimes, \overline{0}, \overline{1})$ be a semiring, and $\leq$ a partial ordering over $A$. We say $K$ is **monotonic** if for all $a, b, c \in A$

$$(a \leq b) \Rightarrow (a \otimes c \leq b \otimes c) \qquad (a \leq b) \Rightarrow (c \otimes a \leq c \otimes b)$$

- **optimal substructure** in dynamic programming

A: $b \otimes c$

B: $b$    C: $c$

A: $b' \otimes c$ $\leq b \otimes c$

B: $b' \leq b$    C: $c$

- idempotent => monotone (from distributivity)

  - (a+b)⊗c = (a⊗c)+(b⊗c);  if a≤b,  (a⊗c)=(a⊗c)+(b⊗c)

  - by def. of comparison,  a⊗c ≤ b⊗c

# DP on Graphs

- optimization problems on graphs
  => generic shortest-path problem

- weighted directed graph G=(V, E) with a function *w* that assigns each edge a weight from a semiring

- compute the best weight of the target vertex *t*

- generic update along edge (u, v)



$$d(v) \oplus = d(u) \otimes w(u,v)$$

$$d(v) \leftarrow d(v) \oplus (d(u) \otimes w(u,v))$$

- how to avoid cyclic updates?
  - only update when d(u) is fixed

# Two Dimensional Survey

| | topological (acyclic) | best-first (superior) |
|---|---|---|
| graphs with semirings (e.g., FSMs) | Viterbi | Dijkstra |
| hypergraphs with weight functions (e.g., CFGs) | Generalized Viterbi | Knuth |

search space

# Viterbi Algorithm for DAGs

1. topological sort

2. visit each vertex v in sorted order and do updates

   - for each incoming edge (u, v) in E

   - use d(u) to update d(v):  $d(v) \oplus = d(u) \otimes w(u, v)$

   - key observation: d(u) is fixed to optimal at this time



   - time complexity: O(V + E )

# Variant 1: forward-update

1. topological sort

2. visit each vertex v in sorted order and do updates

- for each outgoing edge (v, u) in E

- use d(v) to update d(u):  $d(u) \oplus = d(v) \otimes w(v, u)$

- key observation: d(v) is fixed to optimal at this time



- time complexity: O(V + E )

# Examples

# Examples

- [Number of Paths in a DAG]

# Examples

- [Number of Paths in a DAG]

  - just use the counting semiring $(N, +, \times, 0, 1)$

  - note: this is not an optimization problem!

# Examples

- [Number of Paths in a DAG]

  - just use the counting semiring $(N, +, \times, 0, 1)$

  - note: this is not an optimization problem!

- [Longest Path in a DAG]

# Examples

- [Number of Paths in a DAG]

  - just use the counting semiring $(\mathbb{N}, +, \times, 0, 1)$

  - note: this is not an optimization problem!

- [Longest Path in a DAG]

  - just use the semiring $(\mathbb{R} \cup \{-\infty\}, \max, +, -\infty, 0)$

# Examples

- [Number of Paths in a DAG]

  - just use the counting semiring $(\mathbb{N}, +, \times, 0, 1)$

  - note: this is not an optimization problem!

- [Longest Path in a DAG]

  - just use the semiring $(\mathbb{R} \cup \{-\infty\}, \max, +, -\infty, 0)$

- [Part-of-Speech Tagging with a Hidden Markov Model]

# Examples

- [Number of Paths in a DAG]

  - just use the counting semiring $(\mathbb{N}, +, \times, 0, 1)$

  - note: this is not an optimization problem!

- [Longest Path in a DAG]

  - just use the semiring $(\mathbb{R} \cup \{-\infty\}, \max, +, -\infty, 0)$

- [Part-of-Speech Tagging with a Hidden Markov Model]

# Example: Speech Alignment



time complexity: $O(n^2)$

also used in:
edit distance
biological sequence alignment

# Example: Word Alignment



- key difference
  - reorderings in translation!
  - sequence/speech alignment is always monotonic
- complexity under HMM
  - word alignment is $O(n^3)$
    - for every $(i, j)$
      - enumerate all $(i\text{-}1, k)$
  - sequence alignment $O(n^2)$

# Chinese Word Segmentation



下 雨 天 地 面 积 水
xia yu tian di mian ji shui

# Chinese Word Segmentation

民主
**min**-zhu
*people*-*dominate*

"democracy"

下 雨 天 地 面 积 水
xia  yu  tian  di  mian  ji  shui

# Chinese Word Segmentation

民主

min-zhu

*people-dominate*

"democracy"

➡ Google ➡

江泽民 主席

jiang-ze-min  zhu-xi

*... - ... - people  dominate-podium*

"President  Jiang Zemin"

下 雨 天 地 面 积 水
xia  yu  tian  di  mian  ji  shui

# Chinese Word Segmentation

民主

min-zhu

*people*-*dominate*

"democracy"

this was 5 years ago.

now Google is
good at segmentation!

江泽民　主席

jiang-ze-min　zhu-xi

... - ... - *people*　*dominate*-*podium*

"President　Jiang Zemin"

下 雨 天 地 面 积 水

xia  yu  tian  di  mian  ji  shui

# Chinese Word Segmentation

民主

min-zhu

*people*-*dominate*

"democracy"

this was 5 years ago.

now Google is
good at segmentation!

江泽民　主席

jiang-ze-min　zhu-xi

*... - ... - people*　*dominate*-*podium*

"President　Jiang Zemin"

下 雨 天 地 面 积 水
xia yu tian di mian ji shui

# Chinese Word Segmentation

民主

min-zhu

*people-dominate*

"democracy"

Google

this was 5 years ago.

now Google is good at segmentation!

江泽民 主席

jiang-ze-min  zhu-xi

... - ... - *people*  *dominate*-*podium*

"President  Jiang Zemin"

下 雨 天 地 面 积 水
xia  yu  tian  di  mian  ji  shui

graph search

# Phrase-based Decoding

与　沙龙　　举行　了　会谈

*yu Shalong*　*juxing le huitan*

held a talk　　with Sharon

| with | Sharon held | talks |
|------|-------------|-------|

| with Sharon | held a talk |
|-------------|-------------|

*yu　Shalong　juxing　le　huitan*

# Phrase-based Decoding

与 沙龙 举行 了 会谈

*yu Shalong* *juxing le huitan*

held a talk    with Sharon

| with | Sharon held | talks |
|------|-------------|-------|

| with Sharon | held a talk |
|-------------|-------------|

*yu Shalong juxing le huitan*

# Phrase-based Decoding

与　沙龙　　举行　了　会谈

*yu Shalong*　　*juxing le huitan*

held a talk　　with Sharon

| with | Sharon held | talks |
| --- | --- | --- |

| with Sharon | held a talk |
| --- | --- |

*yu　Shalong　juxing　le　huitan*

# Phrase-based Decoding

与　沙龙　　举行　了　会谈

*yu Shalong*　*juxing le huitan*

held a talk　　with Sharon



| with | Sharon held | talks |
|------|-------------|-------|

| with Sharon | held a talk |
|-------------|-------------|

*yu　Shalong*　*juxing　le　huitan*

# Phrase-based Decoding

与　沙龙　　举行　了　会谈

*yu Shalong*　*juxing le huitan*

held a talk　　with Sharon

| with | Sharon held | talks |

| with Sharon | held a talk |

*yu*　*Shalong*　*juxing*　*le*　*huitan*

# Phrase-based Decoding

与　沙龙　　举行　了　会谈

*yu Shalong*　*juxing le huitan*

held a talk　　with Sharon



| with | Sharon held | talks |
|------|-------------|-------|

| with Sharon | held a talk |
|-------------|-------------|

*yu　Shalong　juxing　le　huitan*

# Phrase-based Decoding

与 沙龙 举行 了 会谈

*yu Shalong*  *juxing le huitan*

held a talk  with Sharon

source-side: coverage vector

held a talk

target-side: grow hypotheses
strictly left-to-right

...    ...

held a talk    held a talk with Sharon

...    ...

space: $O(2^n)$, time: $O(2^n n^2)$ -- cf. traveling salesman problem

# Traveling Salesman Problem & MT

- a classical NP-hard problem

  - goal: visit each city once and only once

- exponential-time dynamic programming

  - state: cities visited so far (bit-vector)

  - search in this $O(2^n)$ transformed graph

- MT: each city is a source-language word

  - restrictions in reordering can reduce complexity => distortion limit

  - => syntax-based MT

# Traveling Salesman Problem & MT

- a classical NP-hard problem

  - goal: visit each city once and only once

- exponential-time dynamic programming

  - state: cities visited so far (bit-vector)

  - search in this $O(2^n)$ transformed graph

- MT: each city is a source-language word

  - restrictions in reordering can reduce complexity => distortion limit

  - => syntax-based MT

# Adding a Bigram Model

- "refined" graph: annotated with language model words

- still dynamic programming, just larger search space

# Adding a Bigram Model

- "refined" graph: annotated with language model words

- still dynamic programming, just larger search space

# Adding a Bigram Model

- "refined" graph: annotated with language model words

- still dynamic programming, just larger search space

# Adding a Bigram Model

- "refined" graph: annotated with language model words

- still dynamic programming, just larger search space



$$\text{space: } O(2^n), \quad \text{time: } O(2^n n^2)$$

$$\Rightarrow \quad \text{space: } O(2^n V^{m-1}), \quad \text{time: } O(2^n V^{m-1} n^2)$$

for $m$-gram language models

# Two Dimensional Survey

traversing order

| | topological (acyclic) | best-first (superior) |
|---|---|---|
| graphs with semirings (e.g., FSMs) | Viterbi | Dijkstra |
| hypergraphs with weight functions (e.g., CFGs) | Generalized Viterbi | Knuth |

search space

# Dijkstra Algorithm

$d(u)$ ⬤ ——$w(e)$——▶ ⬤ $d(u) \otimes w(e)$

# Dijkstra Algorithm

- Dijkstra does not require acyclicity

  - instead of topological order, we use best-first order

- but this requires *superiority* of the semiring

  Let $K = (A, \oplus, \otimes, \overline{0}, \overline{1})$ be a semiring, and $\leq$ a partial ordering over $A$. We say $K$ is **superior** if for all $a, b \in A$

  $$a \leq a \otimes b, \qquad b \leq a \otimes b.$$

  - intuition: combination always gets worse

d(u) ⚬ ——$w(e)$——▶ ⚬ d(u) ⊗ *w*(e)

# Dijkstra Algorithm

- Dijkstra does not require acyclicity

  - instead of topological order, we use best-first order

- but this requires *superiority* of the semiring

  Let $K = (A, \oplus, \otimes, \overline{0}, \overline{1})$ be a semiring, and $\leq$ a partial ordering over $A$. We say $K$ is **superior** if for all $a, b \in A$

  $$a \leq a \otimes b, \qquad b \leq a \otimes b.$$

  - intuition: combination always gets worse

- contrast: monotonicity: combination preserves order

  $$(a \leq b) \Rightarrow (a \otimes c \leq b \otimes c)$$

d(u) ◯ ———$w$(e)———▶ ◯ d(u) ⊗ $w$(e)

# Dijkstra Algorithm

- Dijkstra does not require acyclicity

  - instead of topological order, we use best-first order

- but this requires *superiority* of the semiring

  Let $K = (A, \oplus, \otimes, \overline{0}, \overline{1})$ be a semiring, and $\leq$ a partial ordering over $A$. We say $K$ is **superior** if for all $a, b \in A$

$$a \leq a \otimes b, \qquad b \leq a \otimes b.$$

  - intuition: combination always gets worse

- contrast: monotonicity: combination preserves order

$$(a \leq b) \Rightarrow (a \otimes c \leq b \otimes c)$$

$$(\{0, 1\}, \vee, \wedge, 0, 1)$$
$$([0, 1], \max, \times, 0, 1)$$
$$(\mathbb{R}^{+} \cup \{+\infty\}, \mathbf{min}, +, +\infty, 0)$$
$$(\mathbb{R} \cup \{+\infty\}, \mathbf{min}, +, +\infty, 0)$$

d(u) ○ ———$w$(e)——→ ○ d(u) ⊗ $w$(e)

# Dijkstra Algorithm

- Dijkstra does not require acyclicity

  - instead of topological order, we use best-first order

- but this requires *superiority* of the semiring

  Let $K = (A, \oplus, \otimes, \overline{0}, \overline{1})$ be a semiring, and $\leq$ a partial ordering over $A$. We say $K$ is **superior** if for all $a, b \in A$

$$a \leq a \otimes b, \qquad b \leq a \otimes b.$$

  - intuition: combination always gets worse

- contrast: monotonicity: combination preserves order

$$(a \leq b) \Rightarrow (a \otimes c \leq b \otimes c)$$

$(\{0, 1\}, \vee, \wedge, 0, 1)$ ✔

$([0, 1], \max, \times, 0, 1)$ ✔

$(\mathbb{R}^+ \cup \{+\infty\}, \mathbf{min}, +, +\infty, 0)$ ✔

$(\mathbb{R} \cup \{+\infty\}, \mathbf{min}, +, +\infty, 0)$ ✖

d(u) ◯ —————→ ◯ d(u) ⊗ w(e)

$w(e)$

# Dijkstra Algorithm

- keep a cut (S : V - S) where S vertices are fixed

  - maintain a priority queue Q of V - S vertices

- each iteration choose the best vertex v from Q

  - move v to S, and use d(v) to forward-update others

$$d(u) \oplus = d(v) \otimes w(v, u)$$

S          V - S

time complexity:
O((V+E) lgV) (binary heap)
O(V lgV + E) (fib. heap)

# Dijkstra Algorithm

- keep a cut (S : V - S) where S vertices are fixed

  - maintain a priority queue Q of V - S vertices

- each iteration choose the best vertex v from Q

  - move v to S, and use d(v) to forward-update others

$$d(u) \oplus = d(v) \otimes w(v, u)$$

S ...

v

S          V - S

time complexity:
O((V+E) lgV) (binary heap)
O(V lgV + E) (fib. heap)

# Dijkstra Algorithm

- keep a cut (S : V - S) where S vertices are fixed
  - maintain a priority queue Q of V - S vertices
- each iteration choose the best vertex v from Q
  - move v to S, and use d(v) to forward-update others



$$d(u) \oplus = d(v) \otimes w(v, u)$$

time complexity:
O((V+E) lgV) (binary heap)
O(V lgV + E) (fib. heap)

# Viterbi vs. Dijkstra

- structural vs. algebraic constraints

- Dijkstra only applicable to optimization problems

monotonic optimization problems

# Viterbi vs. Dijkstra

- structural vs. algebraic constraints

- Dijkstra only applicable to optimization problems



monotonic optimization problems

acyclic:
Viterbi

# Viterbi vs. Dijkstra

- structural vs. algebraic constraints

- Dijkstra only applicable to optimization problems



monotonic optimization problems

acyclic: Viterbi

superior: Dijkstra

# Viterbi vs. Dijkstra

- structural vs. algebraic constraints

- Dijkstra only applicable to optimization problems



monotonic optimization problems

acyclic: Viterbi

many NLP problems

superior: Dijkstra

# Viterbi vs. Dijkstra

- structural vs. algebraic constraints

- Dijkstra only applicable to optimization problems



monotonic optimization problems

acyclic: Viterbi

many NLP problems

superior: Dijkstra

forward-backward
(Inside semiring)

# Viterbi vs. Dijkstra

- structural vs. algebraic constraints

- Dijkstra only applicable to optimization problems



monotonic optimization problems

many NLP problems

acyclic: Viterbi

superior: Dijkstra

forward-backward (Inside semiring)

non-probabilistic models

# Viterbi vs. Dijkstra

- structural vs. algebraic constraints

- Dijkstra only applicable to optimization problems



monotonic optimization problems

acyclic: Viterbi

many NLP problems

superior: Dijkstra

forward-backward (Inside semiring)

non-probabilistic models

cyclic FSMs/ grammars

# What if both fail?



monotonic optimization problems

acyclic: Viterbi

many NLP problems

superior: Dijkstra

generalized Bellman-Ford
(CLR, 1990; Mohri, 2002)

or, first do strongly-connected components (SCC)
which gives a DAG; use Viterbi globally on this SCC-DAG;
use Bellman-Ford locally within each SCC

# What if both work?



full Dijkstra is slower than Viterbi

$O((V + E) \lg V)$ vs. $O(V + E)$

but it can finish as early as the target vertex is popped

a $(V + E) \lg V$ vs. $V + E$

*Q: how to (magically) reduce a?*

# A* Search: Intuition

- Dijkstra is "blind" about how far the target is

  - may get "trapped" by obstacles

  - can we be more intelligent about the future?

  - idea: prioritize by s-v distance + v-t estimate

# A* Search: Intuition

- Dijkstra is "blind" about how far the target is

  - may get "trapped" by obstacles

  - can we be more intelligent about the future?

  - idea: prioritize by s-v distance + v-t estimate

# A* Search: Intuition

- Dijkstra is "blind" about how far the target is

  - may get "trapped" by obstacles

  - can we be more intelligent about the future?

  - idea: prioritize by s-v distance + v-t estimate

# A* Heuristic



- h(v): the distance from v to target t

  - ĥ(v) must be an optimistic estimate of h(v): ĥ(v) ≤ h(v)

  - Dijkstra is a special case where ĥ(v) = ī   (0 for dist.)

  - now, prioritize the queue by d(v) ⊗ ĥ(v)

- can stop when target gets popped -- why?

  - optimal subpaths should pop earlier than non-optimal

    - d(v) ⊗ ĥ(v) ≤ d(v) ⊗ h(v) ≤ d(t) ≤ non-optimal paths of t

# How to design a heuristic?

- more of an art than science

- basic idea: projection into coarser space

- cluster:    w'(U,V) = min { w(u, v) | u ∈ U, v ∈ V }

- exact cost in coarser graph is estimate of finer graph

# How to design a heuristic?

- more of an art than science

- basic idea: projection into coarser space

- cluster:   $w'(U,V) = \min \{ w(u, v) \mid u \in U, v \in V \}$

- exact cost in coarser graph is estimate of finer graph

# Viterbi or A*?

- A* intuition: $d(t) \otimes \hat{h}(t)$ ranks higher among $d(v) \otimes \hat{h}(v)$

  - can finish early if lucky

  - actually, $d(t) \otimes \hat{h}(t) = d(t) \otimes h(t) = d(t) \otimes \bar{1} = d(t)$

- with the price of maintaining priority queue - $O(\log V)$

- Q: how early? worth the price?

- if the rank is r, then A* is better when $r/V \log V < 1$

d(v) pool

d(v) $\otimes$ $\hat{h}(v)$ pool

d(t)

d(t)

Dijkstra

34

A*

Dynamic Programming

# Viterbi or A*?

- A* intuition: $d(t) \otimes \hat{h}(t)$ ranks higher among $d(v) \otimes \hat{h}(v)$

  - can finish early if lucky

  - actually, $d(t) \otimes \hat{h}(t) = d(t) \otimes h(t) = d(t) \otimes \bar{1} = d(t)$

- with the price of maintaining priority queue - $O(\log V)$

- Q: how early? worth the price?

- if the rank is r,   then A* is better when   $r/V \log V < 1$



d(v) pool

$d(v) \otimes \hat{h}(v)$ pool

$$r < V / \log V$$

d(t)

d(t)

Dijkstra

Liang Huang (Penn)

34

A*

Dynamic Programming

# Two Dimensional Survey

traversing order

| | topological (acyclic) | best-first (superior) |
|---|---|---|
| graphs with semirings (e.g., FSMs) | Viterbi | Dijkstra |
| hypergraphs with weight functions (e.g., CFGs) | Generalized Viterbi | Knuth |

search space

# Two Dimensional Survey

traversing order

| | topological (acyclic) | best-first (superior) |
|---|---|---|
| graphs with semirings (e.g., FSMs) | Viterbi | Dijkstra |
| hypergraphs with weight functions (e.g., CFGs) | Generalized Viterbi | Knuth |

search space

■ For each diff (<= n)
- ■ For each i (<= n)
  - ■ For each rule X → Y Z
    - ■ For each split point k

```
score[X][i][j] = max score[X][i][j],
                     score(X->YZ) *
                     score[Y][i][k] *
                     score[Z][k][j]
```



X

Y    Z

i          k          j

(S, 0, n)

w_0  w_1      ...      w_{n-1}

# Background: CFG and Parsing



- For each diff (<= n)
  - For each i (<= n)
    - For each rule X → Y Z
      - For each split point k

```
score[X][i][j] = max score[X][i][j],
                     score(X->YZ) *
                     score[Y][i][k] *
                     score[Z][k][j]
```

X

Y    Z

i          k          j

(S, 0, n)

w₀  w₁      ...      wₙ₋₁

# Background: CFG and Parsing

- For each diff (<= n)
  - For each i (<= n)
    - For each rule X → Y Z
      - For each split point k
        ```
        score[X][i][j] = max score[X][i][j],
                             score(X->YZ) *
                             score[Y][i][k] *
                             score[Z][k][j]
        ```

- For each diff (<= n)
  - For each i (<= n)
    - For each rule X → Y Z
      - For each split point k

```
score[X][i][j] = max score[X][i][j],
                     score(X->YZ) *
                     score[Y][i][k] *
                     score[Z][k][j]
```



X

Y    Z

i            k            j

(S, 0, n)

$w_0$  $w_1$     ...     $w_{n-1}$

# (Directed) Hypergraphs

- a generalization of graphs

  - edge => hyperedge: several vertices to one vertex

  - $e = (T(e), h(e), f_e)$.    arity $|e| = |T(e)|$

  - a totally-ordered weight set R

    - we borrow the $\oplus$ operator to be the comparison

  - weight function $f_e : R^{|e|}$ to R

    - generalizes the $\otimes$ operator in semirings

simple case:  $f_e(a, b) = a \otimes b \otimes w(e)$

$$d(v) \oplus = f_e(d(u_1), d(u_2))$$

tails

head

# Hypergraphs and Deduction

(B, i, k)    (C, k, j)

$u_1$ : a    $u_2$ : b

$f_e$

v  : a × b × Pr(A → B C)

(A, i, j)

$$\frac{(B, i, k) \qquad (C, k, j)}{(A, i, j)} \quad A \to B\ C$$

(Nederhof, 2003)

# Hypergraphs and Deduction

(B, i, k)        (C, k, j)

$u_1$ : a        $u_2$ : b

$f_e$

v : a × b × Pr(A → B C)

(A, i, j)

$$\frac{(B, i, k) \qquad (C, k, j)}{(A, i, j)} \quad A{\rightarrow}B\ C$$

(Nederhof, 2003)

tails   $u_1$ : a        $u_2$ : b

$f_e$

v : $f_e$ (a,b)

head

$u_1$ : a        $u_2$ : b   antecedents

————————— $f_e$

v : $f_e$ (a,b)

consequent

# Related Formalisms

| hypergraph | AND/OR graph | context-free grammar | deductive system |
|:---:|:---:|:---:|:---:|
| vertex | OR-node | symbol | item |
| source-vertex | leaf OR-node | terminal | axiom |
| target-vertex | root OR-node | start symbol | goal item |
| hyperedge $(\{u_1, u_2\}, v, f)$ | AND-node | production $v \xrightarrow{f} u_1 \ u_2$ | instantiated deduction $$\frac{u_1 : a \quad u_2 : b}{v : f(a, b)}$$ |



OR-node

AND-node

OR-nodes

# Packed Forests

- a compact representation of many parses

  - by sharing common sub-derivations

  - polynomial-space encoding of exponentially large set



$$VP_{1,6}$$
$$e_2 \quad e_1$$
$$VBD_{1,2} \quad NP_{2,6}$$
$$NP_{2,3} \quad PP_{3,6}$$

$$e_1 \quad \frac{VBD_{1,2} \quad NP_{2,3} \quad PP_{3,6}}{VP_{1,6}}$$

$$_0 \; I \; _1 \; saw \; _2 \; him \; _3 \; with \; _4 \; a \; _5 \; mirror \; _6$$

(Klein and Manning, 2001; Huang and Chiang, 2005)

# Packed Forests

- a compact representation of many parses

  - by sharing common sub-derivations

  - polynomial-space encoding of exponentially large set

nodes $\longrightarrow$ $VP_{1,6}$     hyperedges

$e_2$     $e_1$

$VBD_{1,2}$     $NP_{2,6}$

a hypergraph

$NP_{2,3}$     $PP_{3,6}$

$_0$ I $_1$ saw $_2$ him $_3$ with $_4$ a $_5$ mirror $_6$

$$e_1 \quad \frac{VBD_{1,2} \quad NP_{2,3} \quad PP_{3,6}}{VP_{1,6}}$$

(Klein and Manning, 2001; Huang and Chiang, 2005)

# Weight Functions and Semirings



$f_e(a_1, ..., a_k)$

$f_e$

tails

$u_1$

$u_2$

...

$u_k$

$v$

head

# Weight Functions and Semirings



$$f_e(a_1, ..., a_k) = a_1 \otimes ... \otimes a_k \otimes w(e)$$

special case

tails

head

# Weight Functions and Semirings

$d(u)$ ⚫ ⟶ ⚫ $d(u) \otimes w(e)$

$w(e)$

$f_e(a) = a \otimes w(e)$

$d(u)$ ⚫ ⟶ ⚫ $f_e(d(u))$

$f_e$

tails

$u_1$

$u_2$

...

$u_k$

$f_e$ ⟶ v head

$f_e(a_1, ..., a_k) = a_1 \otimes ... \otimes a_k \otimes w(e)$

special case

# Weight Functions and Semirings

$d(u)$ ⬤ ⟶ ⬤ $d(u) \otimes w(e)$

$w(e)$

$f_e(a) = a \otimes w(e)$

$d(u)$ ⬤ ⟶ ⬤ $f_e(d(u))$

$f_e$

semiring-composed

tails

$u_1$
$u_2$
...
$u_k$

$f_e$

$v$  head

$f_e(a_1, ..., a_k) = a_1 \otimes ... \otimes a_k \otimes w(e)$

special case

# Weight Functions and Semirings

$d(u)$ ⬤ ——— $w(e)$ ——→ ⬤ $d(u) \otimes w(e)$

$f_e(a) = a \otimes w(e)$

$d(u)$ ⬤ ——— $f_e$ ——→ ⬤ $f_e(d(u))$

semiring-composed

tails

$u_1$
$u_2$
...
$u_k$

$f_e$ → v head

$f_e(a_1, ..., a_k) = a_1 \otimes ... \otimes a_k \otimes w(e)$

special case

can also extend monotonicity and superiority to general weight functions

# Generalizing Semiring Properties

- monotonicity

  - semiring:  $a \leq b \implies a \times c \leq b \times c$

  - for all weight function f,  for all $a_1 \ldots a_k$,  for all $i$,
    if  $a'_i \leq a_i$  then  $f(a_1 \ldots a'_i \ldots a_k) \leq f(a_1 \ldots a_i \ldots a_k)$

- superiority

  - semiring:  $a \leq a \times b$,  $b \leq a \times b$

  - for all f,  for all $a_1 \ldots a_k$,  for all $i$,    $a_i \leq f(a_1, \ldots, a_k)$

- acyclicity

  - degenerate a hypergraph back into a graph

# Two Dimensional Survey

traversing order

| | topological (acyclic) | best-first (superior) |
|---|---|---|
| graphs with semirings (e.g., FSMs) | Viterbi | Dijkstra |
| hypergraphs with weight functions (e.g., CFGs) | Generalized Viterbi | Knuth |

search space

# Viterbi Algorithm for DAGs

1. topological sort

2. visit each vertex v in sorted order and do updates

   - for each incoming edge (u, v) in E

   - use d(u) to update d(v):

   - key observation: d(u) is fixed to optimal at this time



$$d(v) \oplus = d(u) \otimes w(u, v)$$

   - time complexity: O( V + E )

# Viterbi Algorithm for DA<span style="color:yellow">H</span>s

1. topological sort

2. visit each vertex v in sorted order and do updates

   - for each incoming hyperedge e = ((u₁, .., u|e|), v, f_e)

   - use d(u_i)'s to update d(v)

   - key observation: d(u_i)'s are fixed to optimal at this time



$$d(v) \oplus = f_e(d(u_1), \cdots, d(u_{|e|}))$$

   - time complexity: O( V + E )   (assuming constant arity)

# Example: CKY Parsing

- parsing with CFGs in Chomsky Normal Form (CNF)

- typical instance of the generalized Viterbi for DAHs

- many variants of CKY ~ various topological ordering

(S, 0, n)　　　　　(S, 0, n)

- For each diff (<= n)
  - For each i (<= n)
    - For each rule X → Y Z
      - For each split point k
        `score[X][i][j] = max`

$O(n^3|P|)$

# Example: CKY Parsing

- parsing with CFGs in Chomsky Normal Form (CNF)

- typical instance of the generalized Viterbi for DAHs

- many variants of CKY ~ various topological ordering

(S, 0, n)

(S, 0, n)

bottom-up

$O(n^3|P|)$

- For each diff (<= n)
  - For each i (<= n)
    - For each rule X → Y Z
      - For each split point k
        `score[X][i][j] = max`

# Example: CKY Parsing

- parsing with CFGs in Chomsky Normal Form (CNF)

- typical instance of the generalized Viterbi for DAHs

- many variants of CKY ~ various topological ordering

(S, 0, n)

(S, 0, n)

bottom-up

left-to-right

$O(n^3|P|)$

- For each diff (<= n)
  - For each i (<= n)
    - For each rule X → Y Z
      - For each split point k
        `score[X][i][j] = max`

# Example: CKY Parsing

- parsing with CFGs in Chomsky Normal Form (CNF)

- typical instance of the generalized Viterbi for DAHs

- many variants of CKY ~ various topological ordering

(S, 0, n)

(S, 0, n)

(S, 0, n)

bottom-up

left-to-right

$O(n^3|P|)$

# Example: CKY Parsing

- parsing with CFGs in Chomsky Normal Form (CNF)

- typical instance of the generalized Viterbi for DAHs

- many variants of CKY ~ various topological ordering

(S, 0, n)

(S, 0, n)

(S, 0, n)

bottom-up

left-to-right

right-to-left

$$O(n^3|P|)$$

# Example: Syntax-based MT

- synchronous context-free grammars (SCFGs)

  - context-free grammar in two dimensions

  - generating pairs of strings/trees simultaneously

  - co-indexed nonterminal further rewritten as a unit

$$VP \rightarrow PP^{(1)} VP^{(2)}, \quad VP^{(2)} PP^{(1)}$$
$$VP \rightarrow \textit{juxing le huitan}, \quad \textbf{held a meeting}$$
$$PP \rightarrow \textit{yu Shalong}, \quad \textbf{with Sharon}$$

# Translation as Parsing

- translation with SCFGs => monolingual parsing

- parse the source input with the source projection

  - build the corresponding target sub-strings in parallel

$$\mathbf{VP} \rightarrow \mathbf{PP}^{(1)} \ \mathbf{VP}^{(2)},$$
$$\mathbf{VP} \rightarrow juxing\ le\ huitan,$$
$$\mathbf{PP} \rightarrow yu\ Shalong,$$

VP$_{1,6}$

PP$_{1,3}$     VP$_{3,6}$

*yu  Shalong*     *juxing  le  huitan*

# Translation as Parsing

- translation with SCFGs => monolingual parsing

- parse the source input with the source projection

  - build the corresponding target sub-strings in parallel

$$\text{VP} \rightarrow \text{PP}^{(1)} \text{VP}^{(2)}, \quad \text{VP}^{(2)} \text{PP}^{(1)}$$
$$\text{VP} \rightarrow \textit{juxing le huitan}, \quad \text{held a meeting}$$
$$\text{PP} \rightarrow \textit{yu Shalong}, \quad \text{with Sharon}$$

VP$_{1,6}$

PP$_{1,3}$    VP$_{3,6}$

*yu  Shalong*    *juxing  le  huitan*

# Translation as Parsing

- translation with SCFGs => monolingual parsing

- parse the source input with the source projection

  - build the corresponding target sub-strings in parallel

$$VP \rightarrow PP^{(1)} \; VP^{(2)}, \qquad VP^{(2)} \; PP^{(1)}$$
$$VP \rightarrow juxing \; le \; huitan, \qquad \text{held a meeting}$$
$$PP \rightarrow yu \; Shalong, \qquad \text{with Sharon}$$

held a talk   with Sharon

$$VP_{1,6}$$

with Sharon          held a talk

$$PP_{1,3} \qquad\qquad VP_{3,6}$$

*yu   Shalong*          *juxing   le   huitan*

# Translation as Parsing

- translation with SCFGs => monolingual parsing

- parse the source input with the source projection

  - build the corresponding target sub-strings in parallel

$$\mathbf{VP} \quad \rightarrow \quad \mathbf{PP}^{(1)} \ \mathbf{VP}^{(2)}, \qquad \mathbf{VP}^{(2)} \ \mathbf{PP}^{(1)}$$
$$\mathbf{VP} \quad \rightarrow \quad juxing\ le\ huitan, \qquad \text{held a meeting}$$
$$\mathbf{PP} \quad \rightarrow \quad yu\ Shalong, \qquad \text{with Sharon}$$

held a talk   with Sharon

VP$_{1,6}$

with Sharon          held a talk

PP$_{1,3}$                    VP$_{3,6}$

complexity: same as
CKY parsing -- O(n³)

yu   Shalong          juxing   le   huitan

# Adding a Bigram Model

# Adding a Bigram Model

# Adding a Bigram Model



complexity: $O(n^3 V^{4(m-1)})$

# Two Dimensional Survey

traversing order

| | topological (acyclic) | best-first (superior) |
|---|---|---|
| graphs with semirings (e.g., FSMs) | Viterbi | Dijkstra |
| hypergraphs with weight functions (e.g., CFGs) | Generalized Viterbi | Knuth |

search space

# Viterbi Algorithm for DAHs

1. topological sort

2. visit each vertex v in sorted order and do updates

   - for each incoming hyperedge e = (($u_1$, .., $u_{|e|}$), v, $f_e$)

   - use $d(u_i)$'s to update d(v)

   - key observation: $d(u_i)$'s are fixed to optimal at this time

$$d(v) \oplus= f_e(d(u_1), \cdots, d(u_{|e|}))$$

   - time complexity: O( V + E )   (assuming constant arity)

# Forward Variant for DA*H*s

1. topological sort

2. visit each vertex v in sorted order and do updates

   - for each outgoing hyperedge $e = ((u_1, .., u_{|e|}), h(e), f_e)$

   - if $d(u_i)$'s have all been fixed to optimal

     - use $d(u_i)$'s to update $d(h(e))$

$v = u_i$



   - time complexity: $O(V + E)$

# Forward Variant for DAHs

1. topological sort

2. visit each vertex v in sorted order and do updates

   - for each outgoing hyperedge e = $((u_1, .., u_{|e|}), h(e), f_e)$

   - if $d(u_i)$'s have all been fixed to optimal

     - use $d(u_i)$'s to update $d(h(e))$



$v = u_i$

- time complexity: O( V + E )

# Forward Variant for DAHs

1. topological sort

2. visit each vertex v in sorted order and do updates

  - for each outgoing hyperedge $e = ((u_1, .., u_{|e|}), h(e), f_e)$

  - if $d(u_i)$'s have all been fixed to optimal

    - use $d(u_i)$'s to update $d(h(e))$

$v = u_i$

Q: *how to avoid repeated checking?*
maintain a counter r[e] for each e:
  how many tails yet to be fixed?
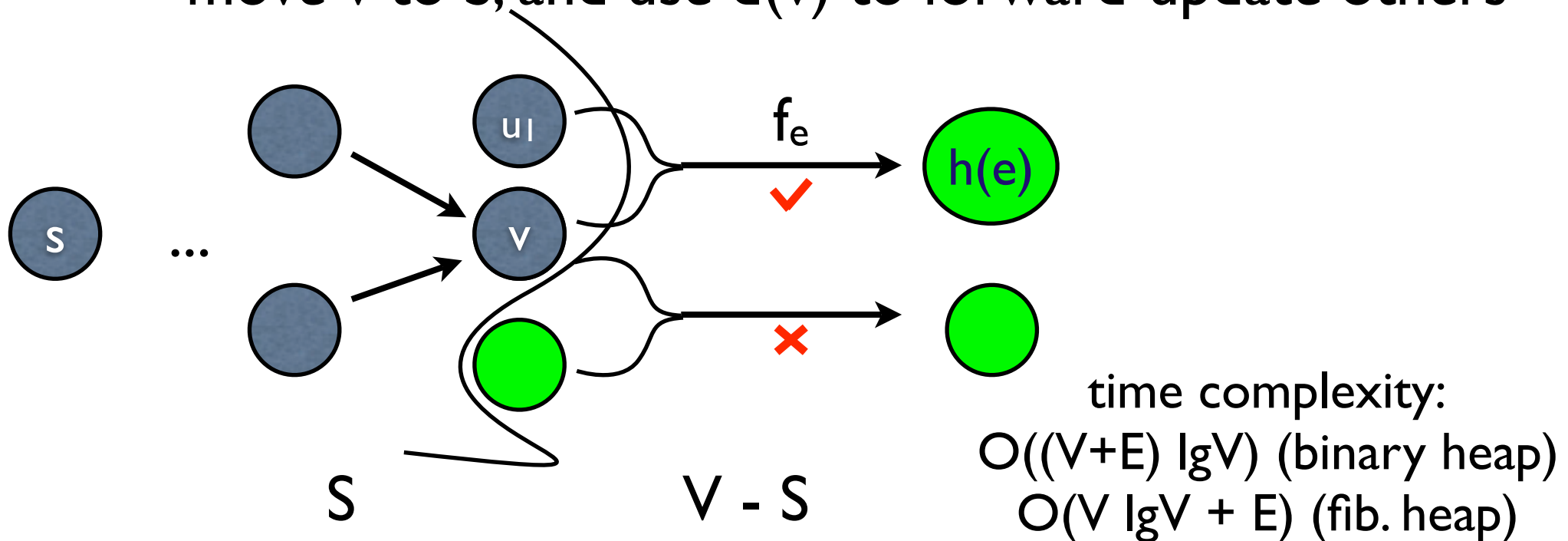fire this hyperedge only if r[e]=0

- time complexity: O( V + E )

# Dijkstra Algorithm

- keep a cut (S : V - S) where S vertices are fixed

  - maintain a priority queue Q of V - S vertices

- each iteration choose the best vertex v from Q
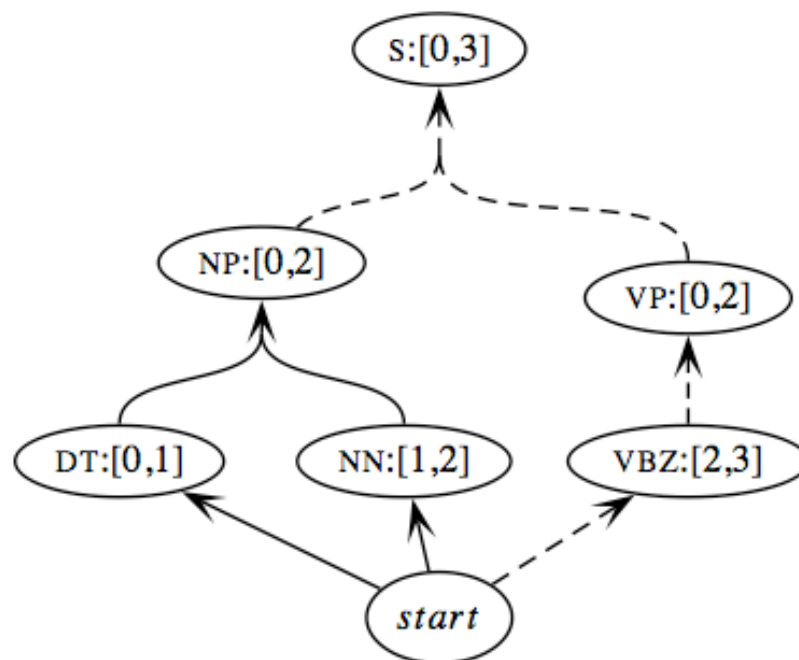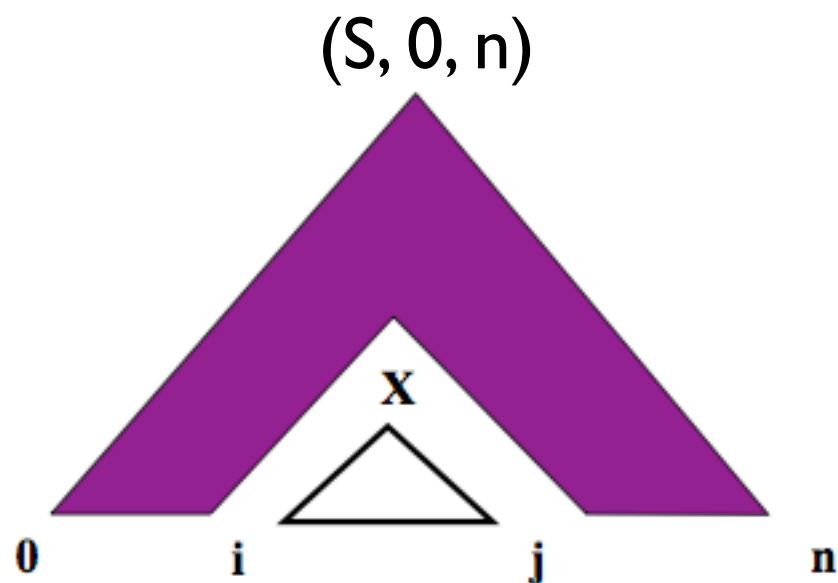
  - move v to S, and use d(v) to forward-update others

$$d(u) \oplus = d(v) \otimes w(v, u)$$

S        ...        v

S                V - S

time complexity:
O((V+E) lgV) (binary heap)
O(V lgV + E) (fib. heap)

# Dijkstra Algorithm

- keep a cut (S : V - S) where S vertices are fixed

  - maintain a priority queue Q of V - S vertices

- each iteration choose the best vertex v from Q
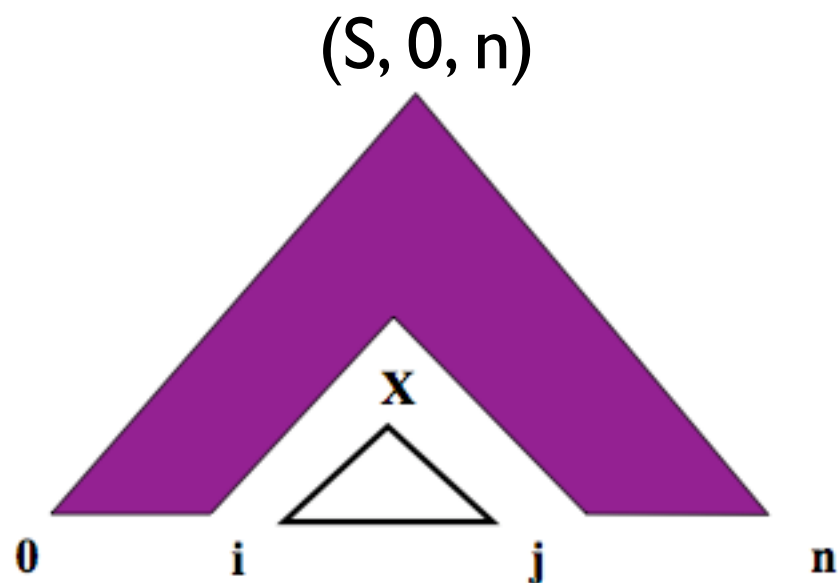
  - move v to S, and use d(v) to forward-update others

$$d(u) \oplus = d(v) \otimes w(v, u)$$



S          V - S

time complexity:
O((V+E) lgV) (binary heap)
O(V lgV + E) (fib. heap)

# Dijkstra Algorithm

- keep a cut (S : V - S) where S vertices are fixed

  - maintain a priority queue Q of V - S vertices

- each iteration choose the best vertex v from Q

  - move v to S, and use d(v) to forward-update others



$$d(u) \oplus = d(v) \otimes w(v, u)$$

w(v, u)

S      ...      v      S      V - S

time complexity:
O((V+E) lgV) (binary heap)
O(V lgV + E) (fib. heap)

# Knuth (1977) Algorithm

- keep a cut (S : V - S) where S vertices are fixed

  - maintain a priority queue Q of V - S vertices

- each iteration choose the best vertex v from Q

  - move v to S, and use d(v) to forward-update others



time complexity:
O((V+E) lgV) (binary heap)
O(V lgV + E) (fib. heap)

# Knuth (1977) Algorithm

- keep a cut (S :V - S) where S vertices are fixed

  - maintain a priority queue Q of V - S vertices

- each iteration choose the best vertex v from Q

  - move v to S, and use d(v) to forward-update others



time complexity:
O((V+E) lgV) (binary heap)
O(V lgV + E) (fib. heap)

# Knuth (1977) Algorithm

- keep a cut (S : V - S) where S vertices are fixed

  - maintain a priority queue Q of V - S vertices

- each iteration choose the best vertex v from Q

  - move v to S, and use d(v) to forward-update others



$f_e$

h(e)

S          V - S

time complexity:
$O((V+E) \lg V)$ (binary heap)
$O(V \lg V + E)$ (fib. heap)

# Knuth (1977) Algorithm

- keep a cut (S : V - S) where S vertices are fixed

  - maintain a priority queue Q of V - S vertices

- each iteration choose the best vertex v from Q

  - move v to S, and use d(v) to forward-update others



$f_e$

✔

✘

h(e)

S

V - S

time complexity:
O((V+E) lgV) (binary heap)
O(V lgV + E) (fib. heap)

# Example: Best-First/A* Parsing

- Knuth for parsing: best-first (Caraballo & Charniak, 1998)

- further speed-up: use A* heuristics

  - showed significant speed up with carefully designed heuristic functions (Klein and Manning, 2003)

  - heuristic function: an estimate of outside cost



(S, 0, n)

# Example: Best-First/A* Parsing

- Knuth for parsing:  best-first  <span style="color:gray">(Caraballo & Charniak, 1998)</span>

- further speed-up: use A* heuristics

  - showed significant speed up with carefully designed heuristic functions <span style="color:gray">(Klein and Manning, 2003)</span>

  - heuristic function: an estimate of outside cost

$(S, 0, n)$

# Outside Cost in Hypergraph

- outside cost: yet to pay to reach goal
- let's only consider semiring-composed case
  - and only acyclic hypergraphs
- after computing d(v) for all v from bottom-up
  - backwards Viterbi from top-down (outside-in)

$h(S_{0,n}) = \bar{1}$

$h(v) \oplus= h(u) \otimes w(e) \otimes d(v')$

# Outside Cost in Hypergraph

- outside cost: yet to pay to reach goal

- let's only consider semiring-composed case

  - and only acyclic hypergraphs

- after computing d(v) for all v from bottom-up

  - backwards Viterbi from top-down (outside-in)

$$h(S_{0,n}) = \bar{1}$$
$$h(v) \oplus= h(u) \otimes w(e) \otimes d(v')$$
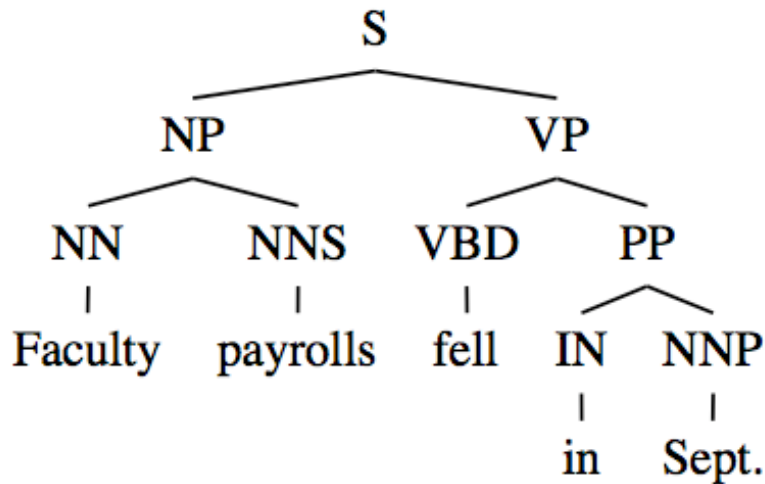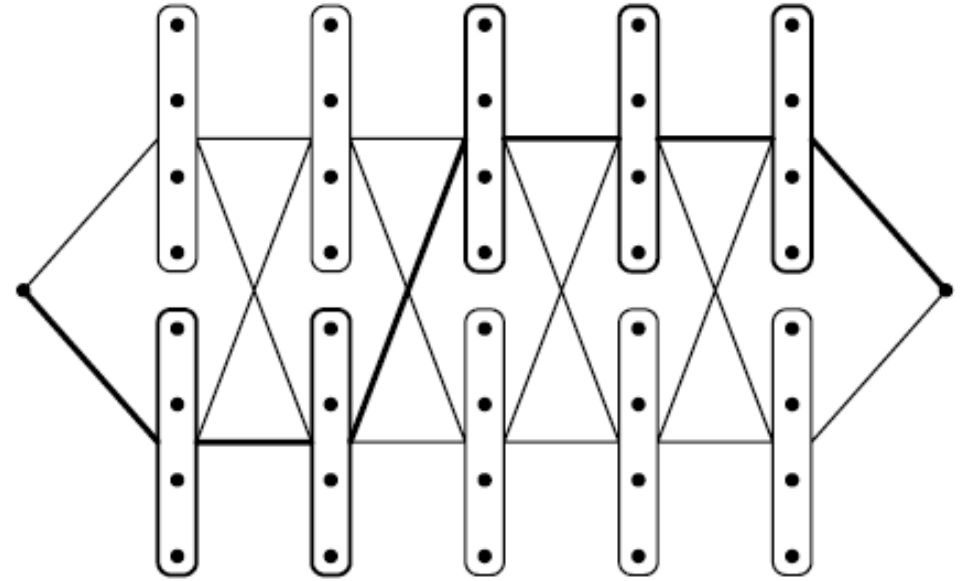
$Q$: $d(v) \otimes h(v) = ?$

# Projection-based Heuristics

- how to guess? project onto a coarser-grained space

- and parse with the coarser grammar

  - outside cost of of the coarser item as heuristics
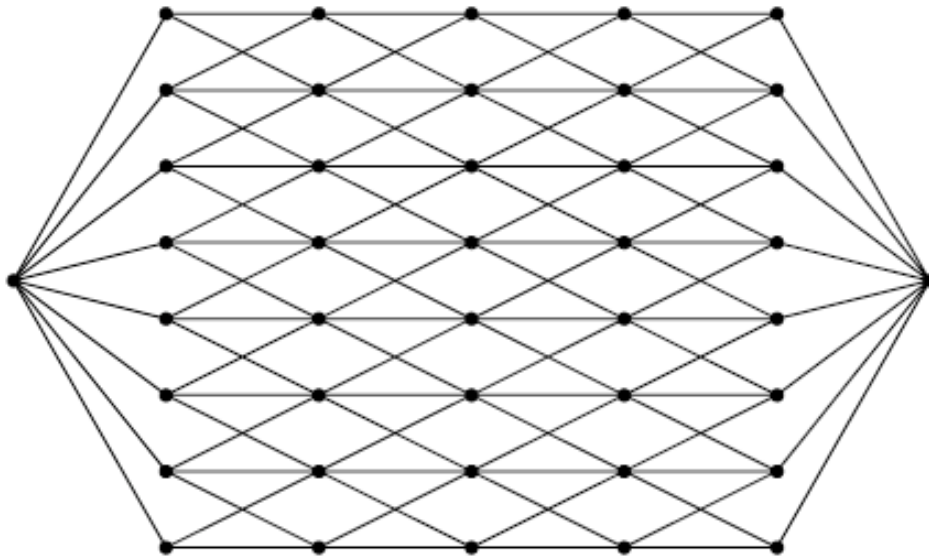


(Klein and Manning, 2003)

# Projection-based Heuristics

- how to guess? project onto a coarser-grained space

- and parse with the coarser grammar

  - outside cost of of the coarser item as heuristics



(Klein and Manning, 2003)

# Projection-based Heuristics

- how to guess? project onto a coarser-grained space

- and parse with the coarser grammar

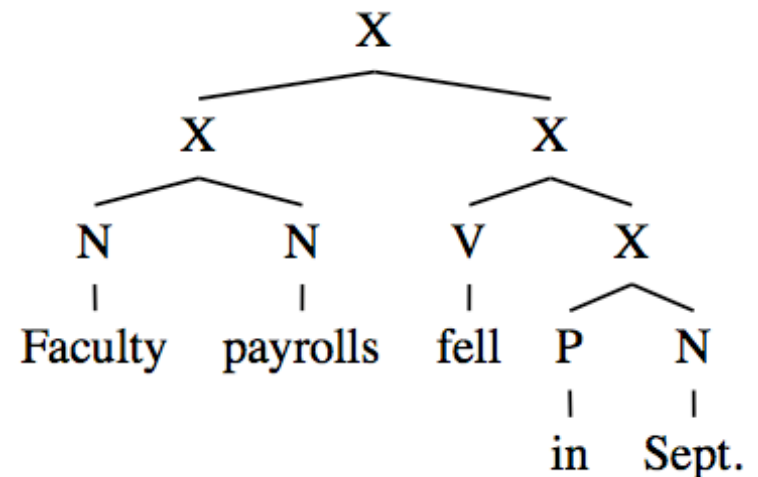  - outside cost of of the coarser item as heuristics
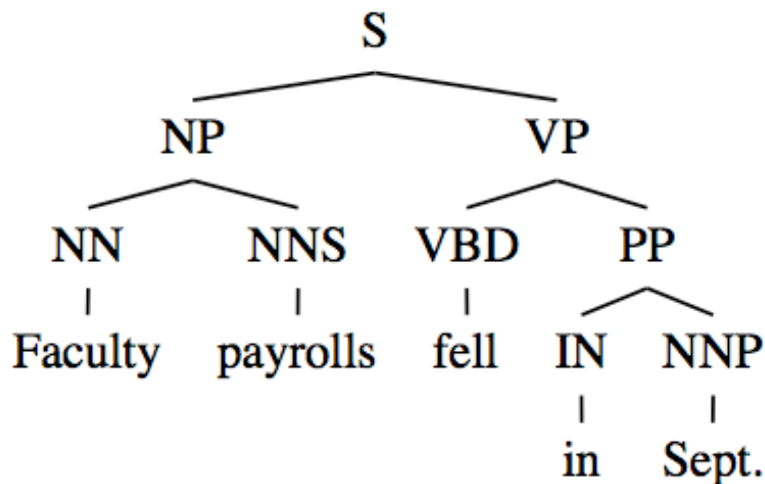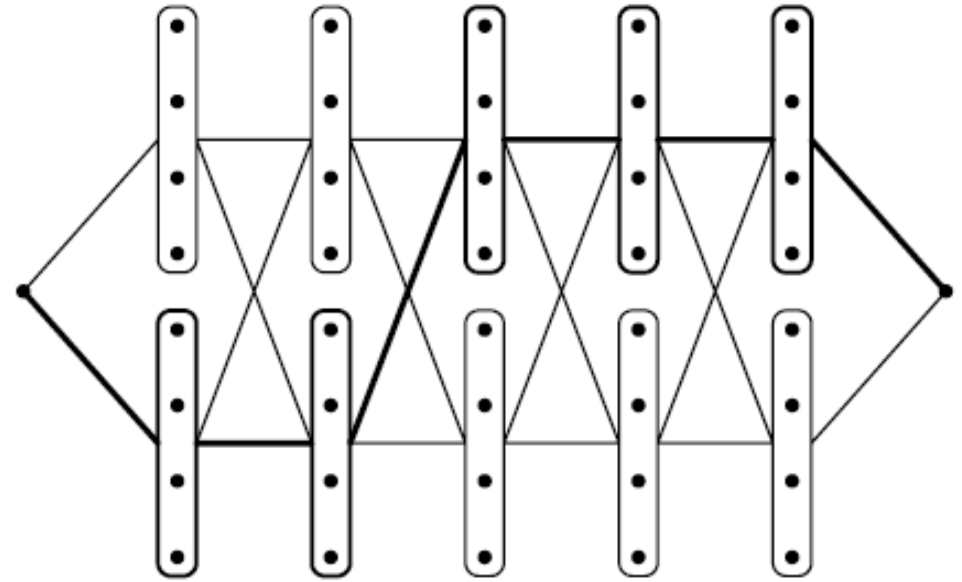


(Klein and Manning, 2003)

# Projection-based Heuristics

- how to guess? project onto a coarser-grained space

- and parse with the coarser grammar
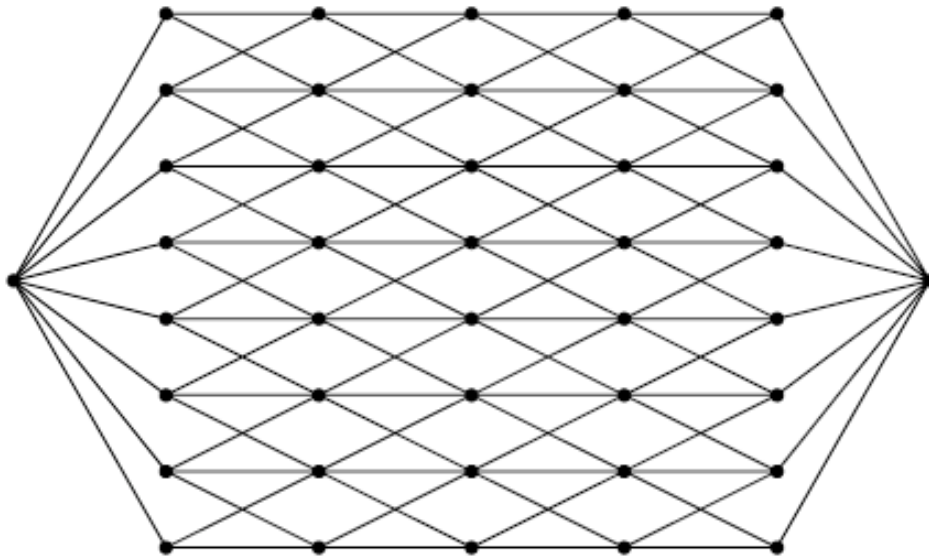
  - outside cost of of the coarser item as heuristics



(Klein and Manning, 2003)

# Projection-based Heuristics

- how to guess? project onto a coarser-grained space

- and parse with the coarser grammar

  - outside cost of of the coarser item as heuristics



(Klein and Manning, 2003)

# Projection-based Heuristics

- how to guess? project onto a coarser-grained space

- and parse with the coarser grammar

  - outside cost of of the coarser item as heuristics



$$\hat{h}\,(\text{VBD}_{2,3}) = h'\,(V_{2,3}) \quad \text{(Klein and Manning, 2003)}$$
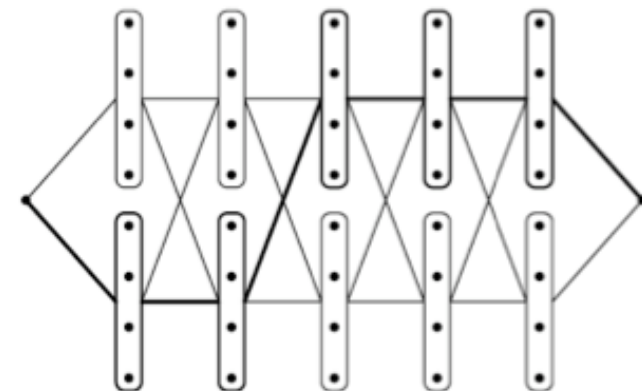
# Analogy with Graphs

# Analogy with Graphs

# More on Coarse-to-Fine

- multilevel coarse-to-fine A*

  - heuristic = exact outside cost in previous stage

  - $\hat{h}_i(v) = h_{i-1}(\text{proj}_{i-1}(v))$

  - VBD>V>X.  $\hat{h}_i(VBD_{1,5}) = h_{i-1}(V_{1,5}); \hat{h}_{i-1}(V_{1,5}) = h_{i-2}(X_{1,5})$

- multilevel coarse-to-fine Viterbi w/ beam-search

  - Viterbi + beam pruning in each stage

  - prune according to merit:  $d(v) \otimes h(v) \oslash d(TOP)$

  - hard to derive a provably correct threshold

  - in practice: use a preset threshold (but works well!)
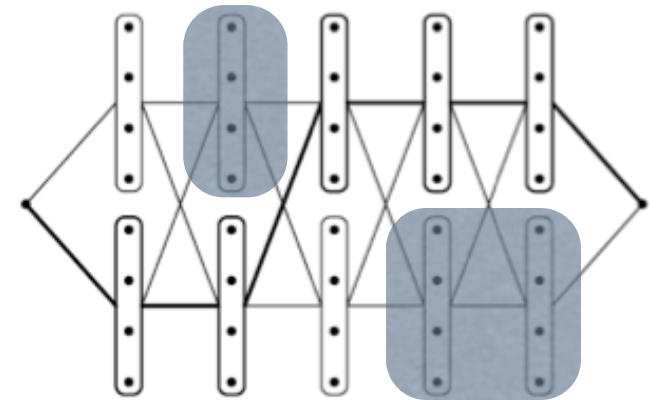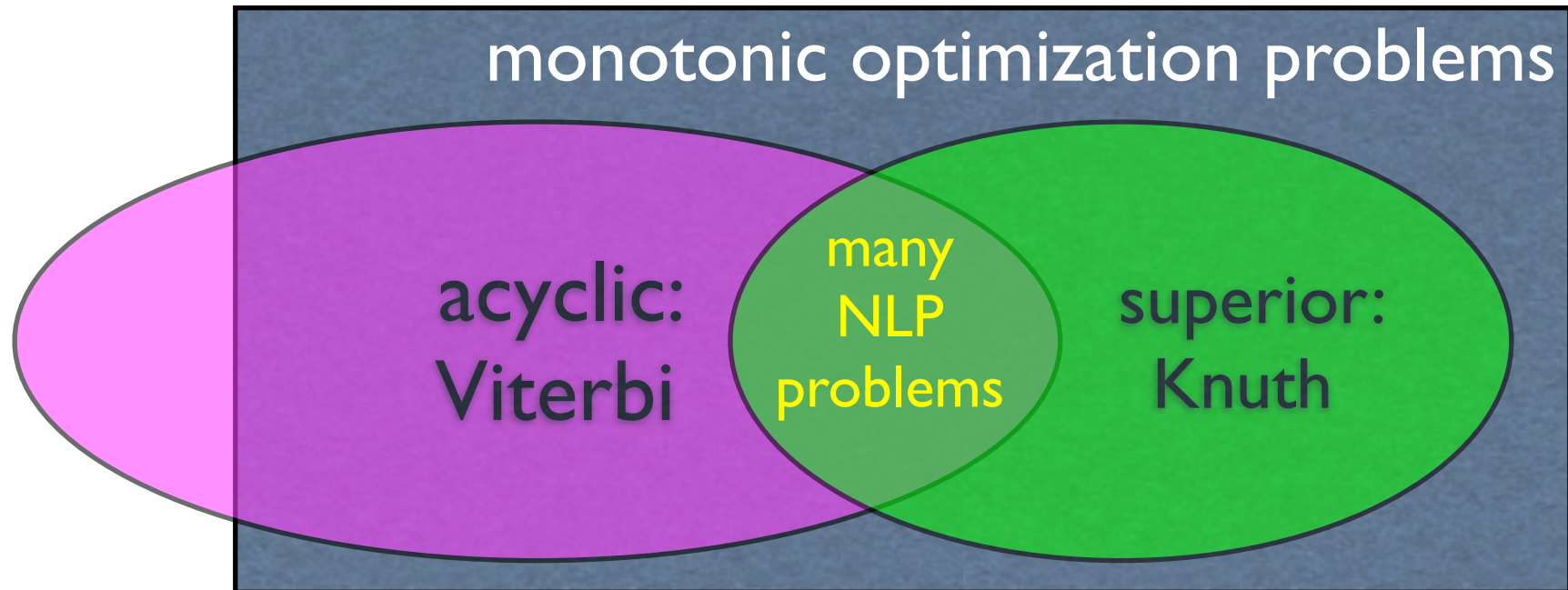
# More on Coarse-to-Fine

- multilevel coarse-to-fine A*

  - heuristic = exact outside cost in 

  - $\hat{h}_i(v) = h_{i-1}(\text{proj}_{i-1}(v))$

  - VBD>V>X.  $\hat{h}_i(VBD_{1,5}) = h_{i-1}(V_{1,5})$; $\hat{h}_{i-1}(V_{1,5}) = h_{i-2}(X_{1,5})$

- multilevel coarse-to-fine Viterbi w/ beam-search

  - Viterbi + beam pruning in each stage

  - prune according to merit:  $d(v) \otimes h(v) \oslash d(TOP)$

  - hard to derive a provably correct threshold

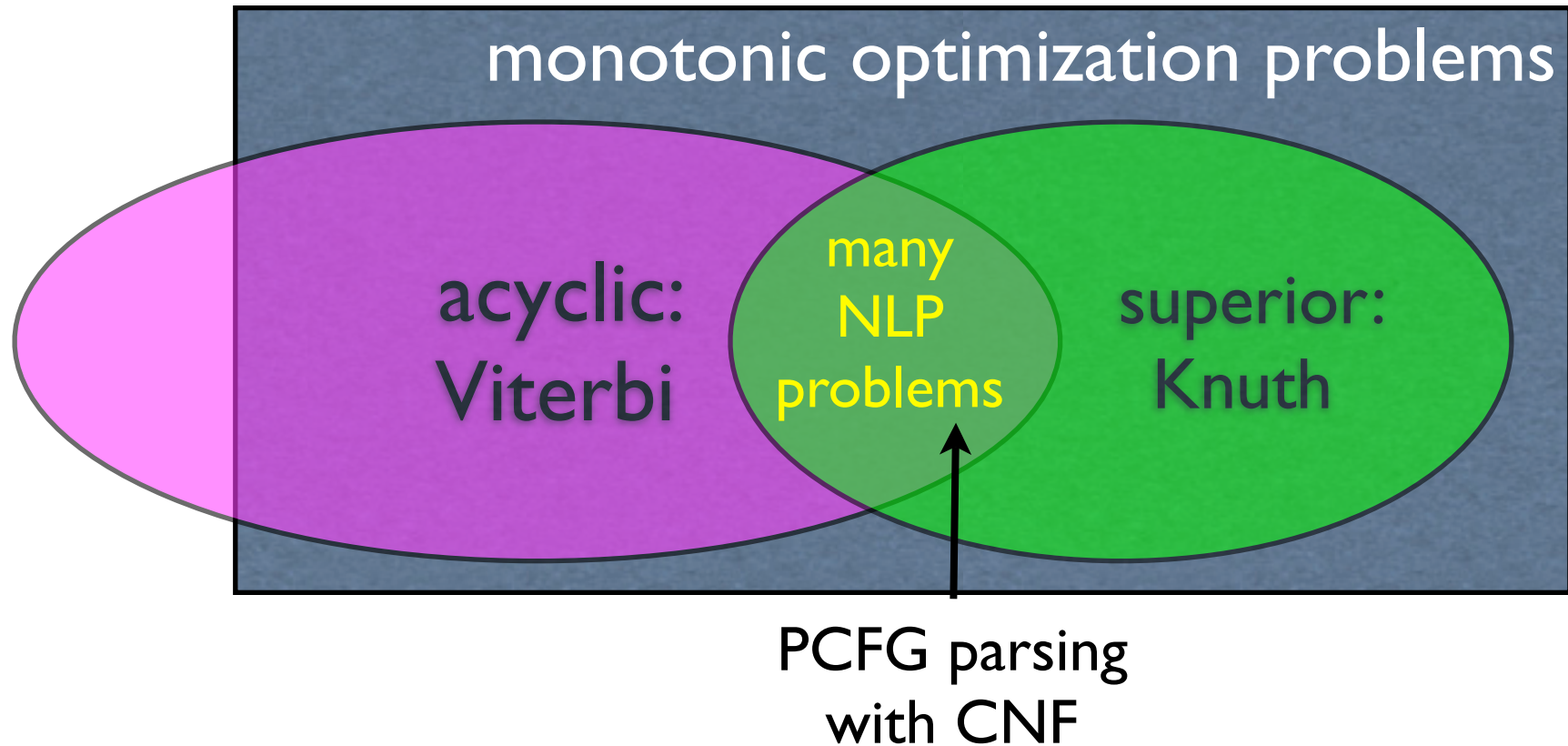  - in practice: use a preset threshold (but works well!)

# More on Coarse-to-Fine

- multilevel coarse-to-fine A*

  - heuristic = exact outside cost in

  - $\hat{h}_i(v) = h_{i-1}(\text{proj}_{i-1}(v))$

  - VBD>V>X.  $\hat{h}_i(VBD_{1,5}) = h_{i-1}(V_{1,5}); \hat{h}_{i-1}(V_{1,5}) = h_{i-2}(X_{1,5})$

- multilevel coarse-to-fine Viterbi w/ beam-search

  - Viterbi + beam pruning in each stage

  - prune according to merit:  $d(v) \otimes h(v) \oslash d(\text{TOP})$

  - hard to derive a provably correct threshold

  - in practice: use a preset threshold (but works well!)

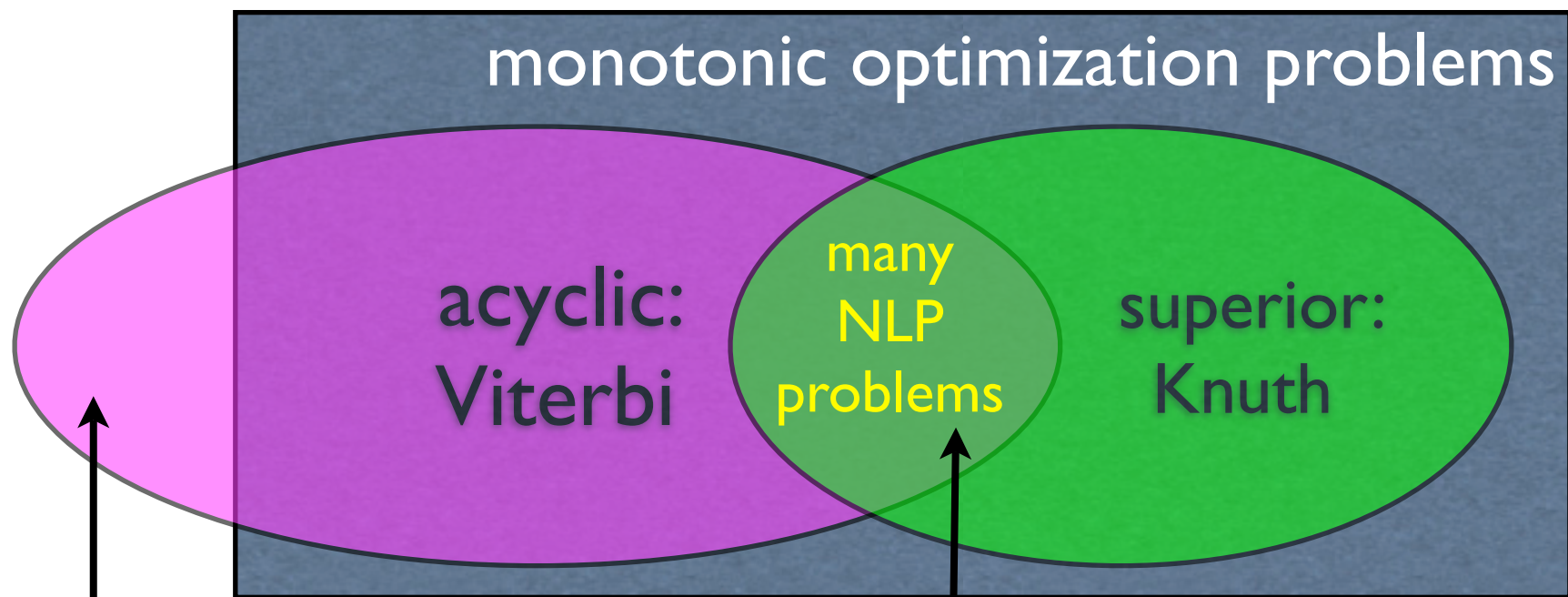# Same Picture Again



monotonic optimization problems

acyclic: Viterbi

many NLP problems

superior: Knuth

# Same Picture Again

# Same Picture Again

monotonic optimization problems
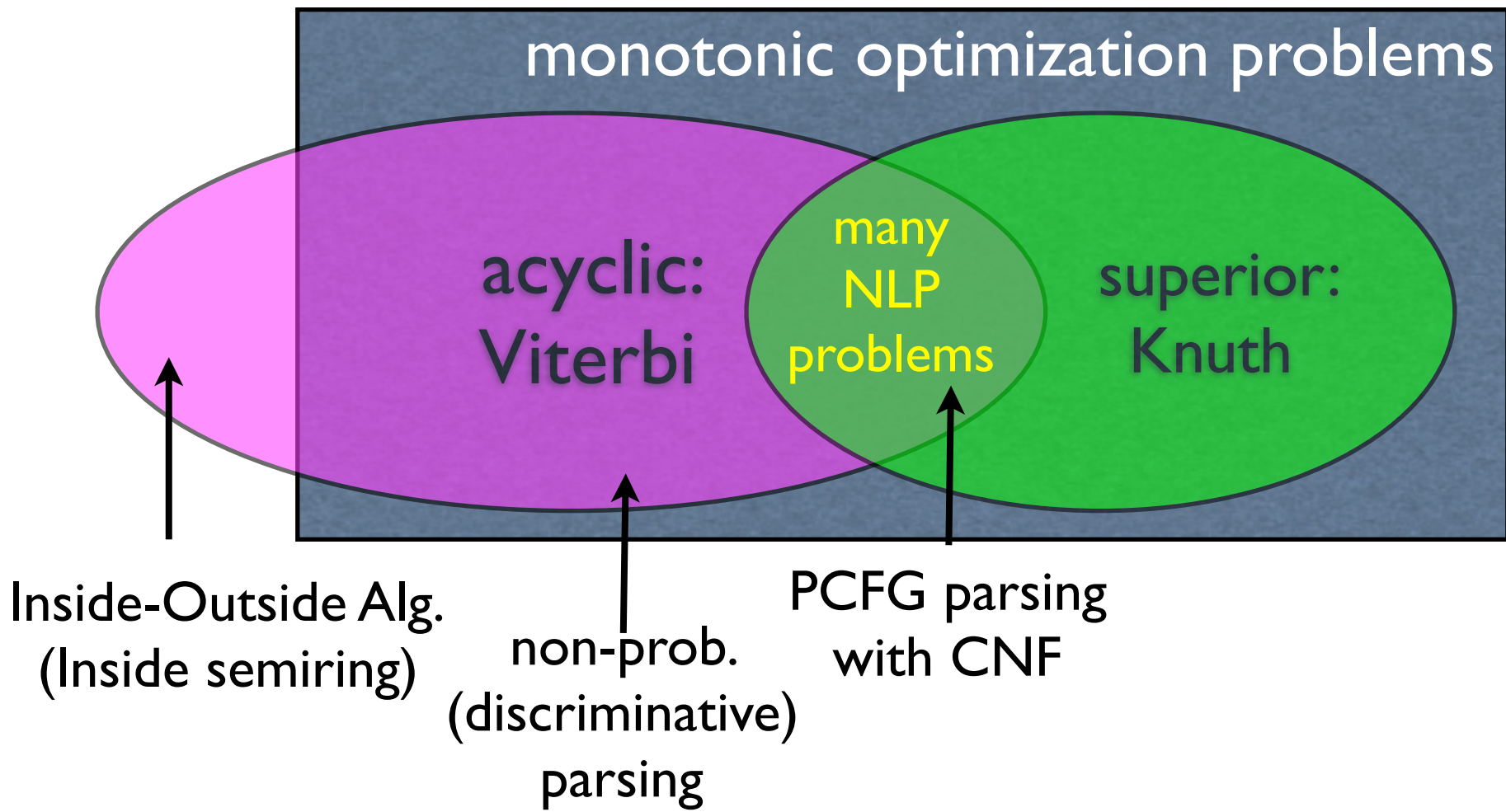
many NLP problems

acyclic: Viterbi

superior: Knuth

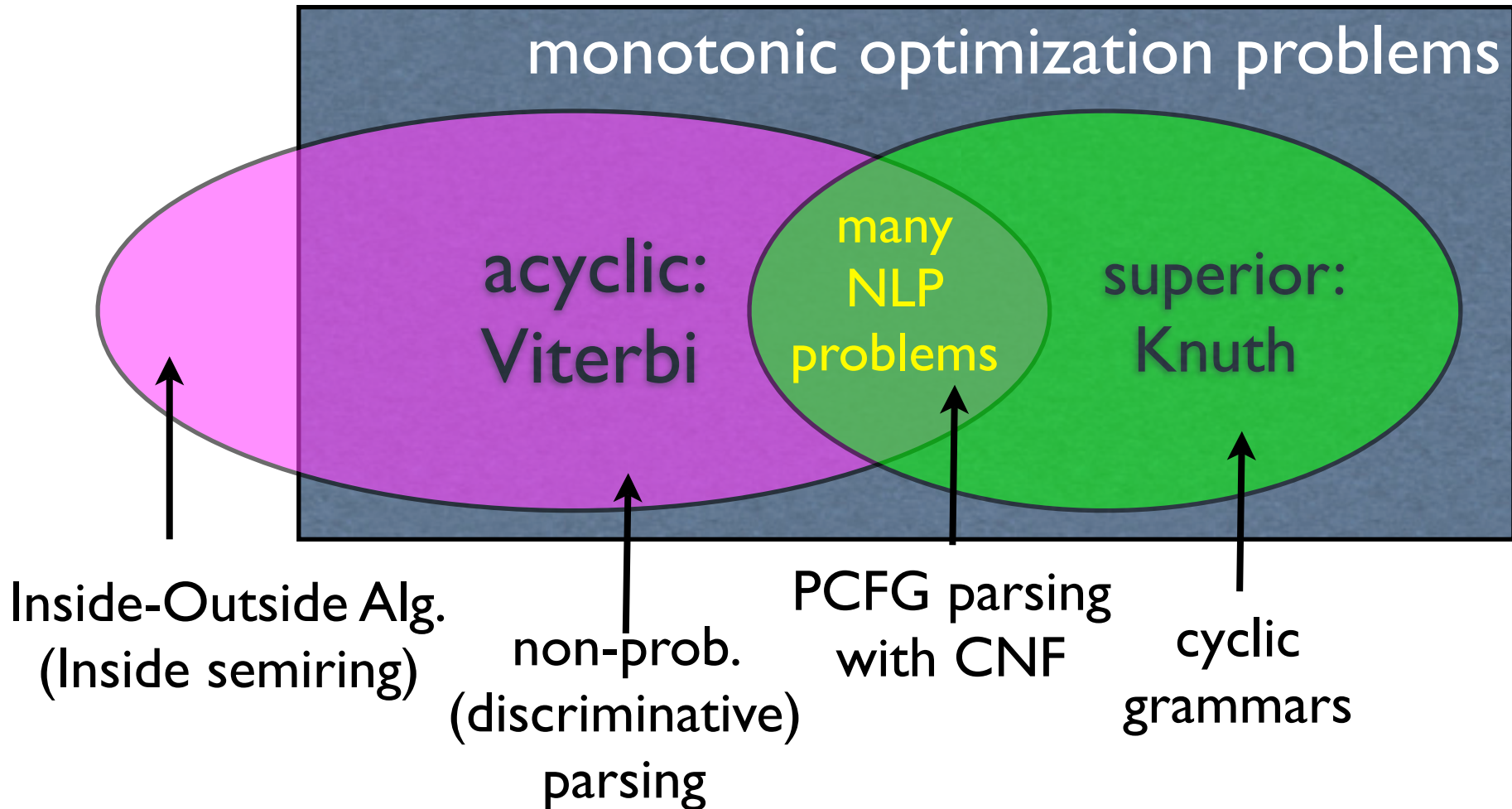Inside-Outside Alg.
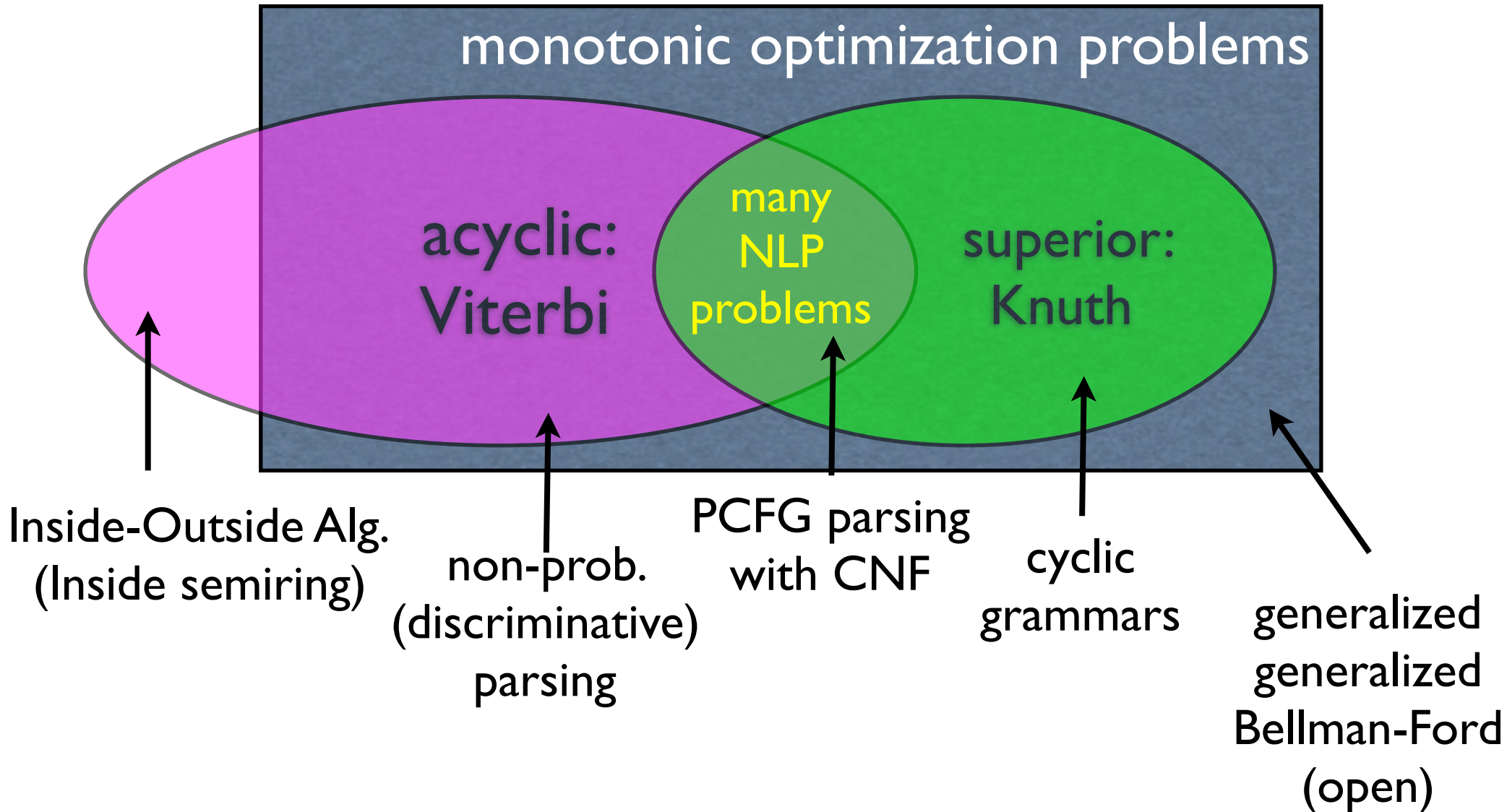(Inside semiring)

PCFG parsing
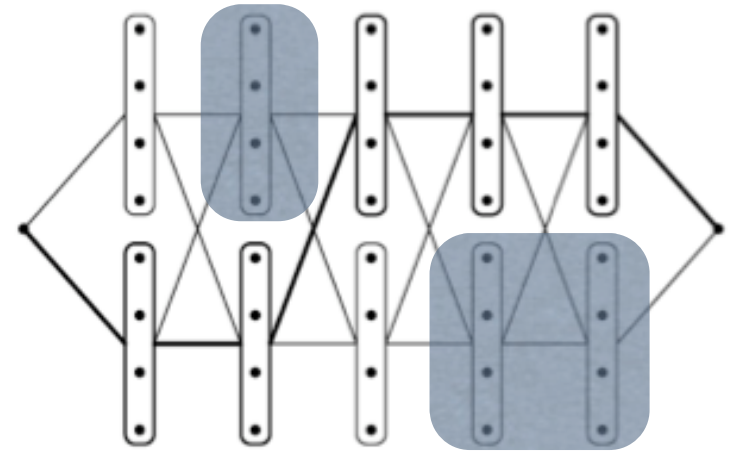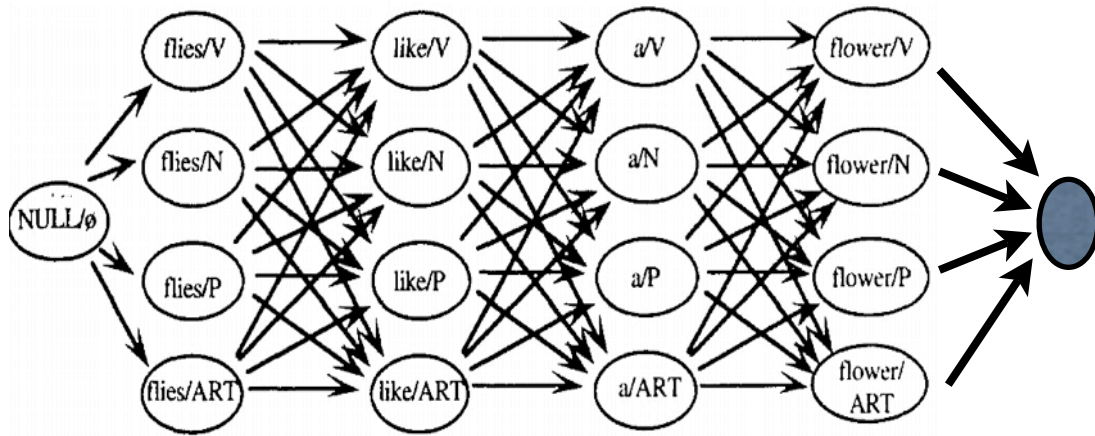with CNF
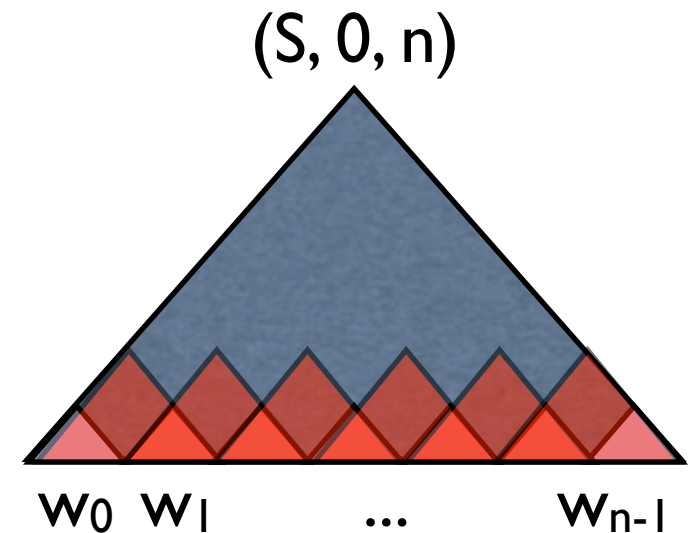
# Same Picture Again

# Same Picture Again

# Same Picture Again

# Take Home Message

- Dynamic Programming is cool, easy, and universal!

- two frameworks and two types of algorithms

  - monotonicity; acyclicity and/or superiority

  - topological (Viterbi) vs. best-first style (Dijkstra/Knuth/A*)

    - when to choose which: A* can finish early if lucky

  - graph (lattice) vs. hypergraph (forest)

    - incremental, finite-state vs. branching, context-free

- covered many typical NLP applications

  - a better understanding of theory helps in practice

# Thanks!



Questions?
Comments?

$(S, 0, n)$

$w_0$ $w_1$ ... $w_{n-1}$

final slides will be available on my website.