

# Dual Decomposition for Parsing with Non-Projective Head Automata

Terry Koo, Alexander M. Rush, Michael Collins,  
David Sontag, and Tommi Jaakkola

# The Cost of Model Complexity

We are always looking for better ways to model natural language.

Tradeoff: Richer models  $\Rightarrow$  Harder decoding

Added complexity is both computational and implementational.

# The Cost of Model Complexity

We are always looking for better ways to model natural language.

Tradeoff: Richer models  $\Rightarrow$  Harder decoding

Added complexity is both computational and implementational.

Tasks with challenging decoding problems:

- ▶ Speech Recognition
- ▶ Sequence Modeling (e.g. extensions to HMM/CRF)
- ▶ Parsing
- ▶ Machine Translation

# The Cost of Model Complexity

We are always looking for better ways to model natural language.

Tradeoff: Richer models  $\Rightarrow$  Harder decoding

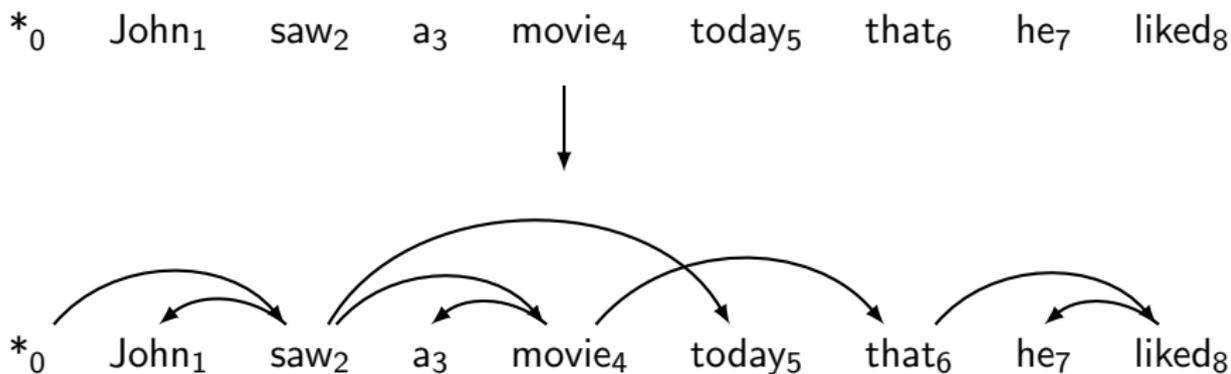
Added complexity is both computational and implementational.

Tasks with challenging decoding problems:

- ▶ Speech Recognition
- ▶ Sequence Modeling (e.g. extensions to HMM/CRF)
- ▶ Parsing
- ▶ Machine Translation

$$y^* = \arg \max_y f(y) \quad \text{Decoding}$$

# Non-Projective Dependency Parsing



Important problem in many languages.

Problem is **NP-Hard** for all but the simplest models.

# Dual Decomposition

A classical technique for constructing decoding algorithms.

Solve complicated models

$$y^* = \arg \max_y f(y)$$

by decomposing into smaller problems.

Upshot: Can utilize a toolbox of combinatorial algorithms.

- ▶ Dynamic programming
- ▶ Minimum spanning tree
- ▶ Shortest path
- ▶ Min-Cut
- ▶ ...

# A Dual Decomposition Algorithm for Non-Projective Dependency Parsing

**Simple** - Uses basic combinatorial algorithms

**Efficient** - Faster than previously proposed algorithms

**Strong Guarantees** - Gives a certificate of optimality when exact

Solves 98% of examples exactly, even though the problem is NP-Hard

**Widely Applicable** - Similar techniques extend to other problems

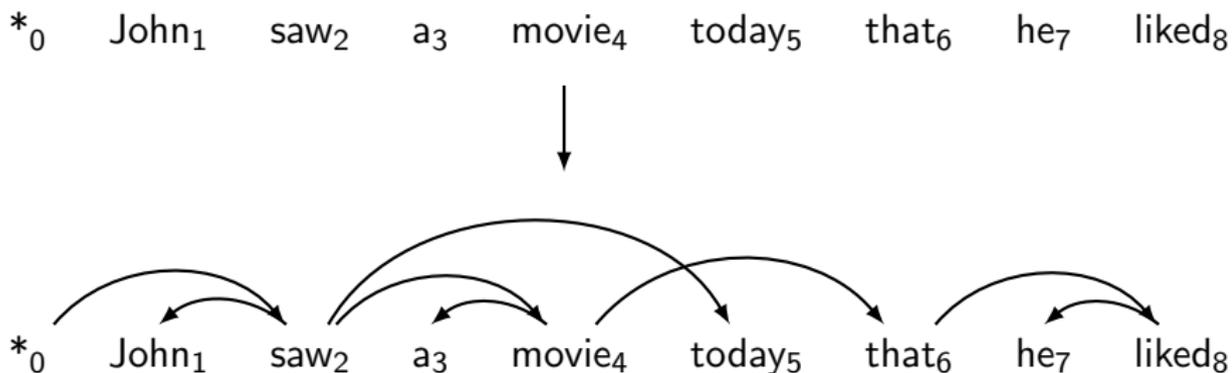
Roadmap

Algorithm

Experiments

Derivation

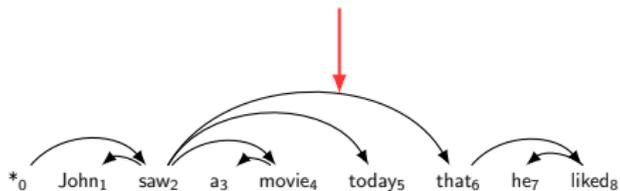
# Non-Projective Dependency Parsing



- ▶ Starts at the root symbol \*
- ▶ Each word has a exactly one parent word
- ▶ Produces a tree structure (no cycles)
- ▶ Dependencies can cross

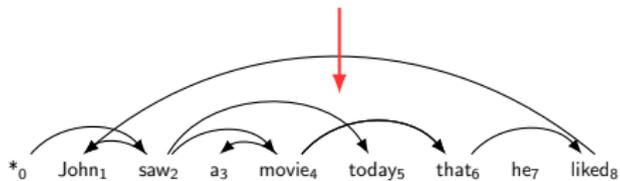
# Algorithm Outline

\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>



Arc-Factored Model

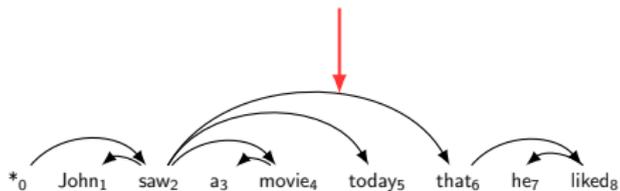
\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>



Sibling Model

# Algorithm Outline

\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>

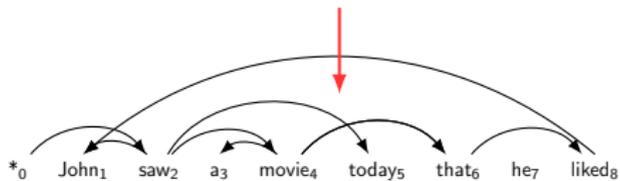


Arc-Factored Model

+

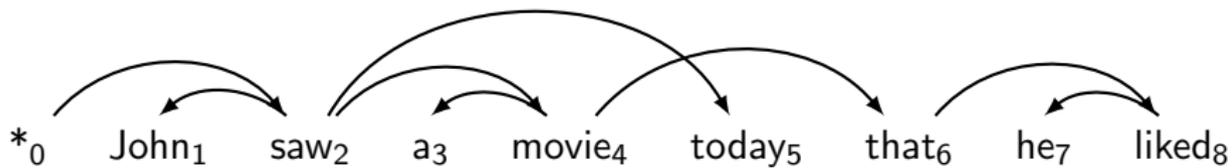
Dual Decomposition

\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>



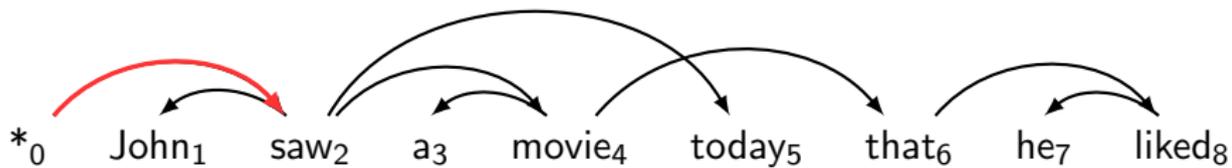
Sibling Model

## Arc-Factored



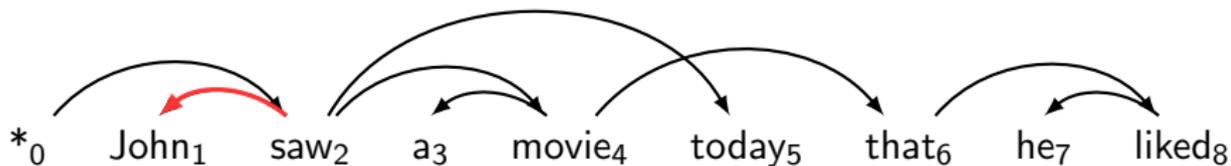
$f(y) =$

## Arc-Factored



$$f(y) = \text{score}(\text{head} = *_0, \text{mod} = \text{saw}_2)$$

## Arc-Factored



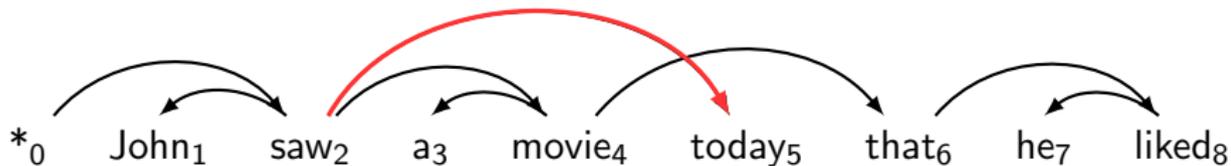
$$f(y) = \text{score}(\text{head} = *_{0}, \text{mod} = \text{saw}_{2}) + \text{score}(\text{saw}_{2}, \text{John}_{1})$$

## Arc-Factored



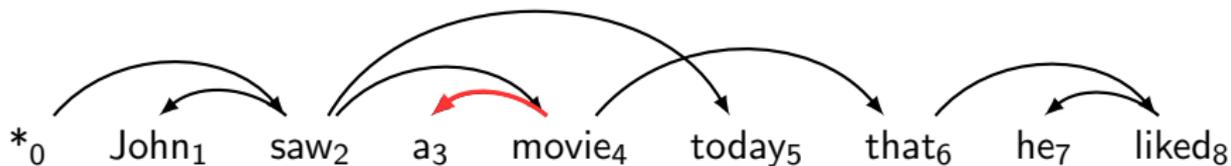
$$f(y) = \text{score}(\text{head} = *_{0}, \text{mod} = \text{saw}_{2}) + \text{score}(\text{saw}_{2}, \text{John}_{1}) \\ + \text{score}(\text{saw}_{2}, \text{movie}_{4})$$

## Arc-Factored



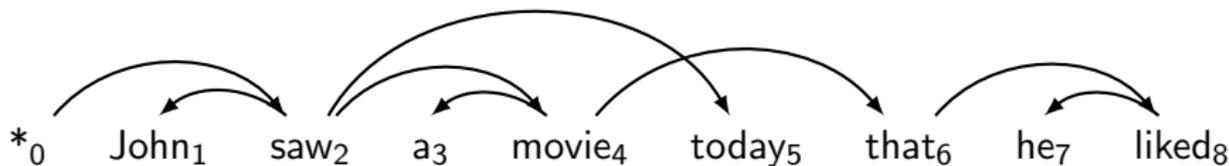
$$f(y) = \text{score}(\text{head} = *_{0}, \text{mod} = \text{saw}_{2}) + \text{score}(\text{saw}_{2}, \text{John}_{1}) \\ + \text{score}(\text{saw}_{2}, \text{movie}_{4}) + \text{score}(\text{saw}_{2}, \text{today}_{5})$$

## Arc-Factored



$$\begin{aligned} f(y) = & \text{score}(\text{head} = *_0, \text{mod} = \text{saw}_2) + \text{score}(\text{saw}_2, \text{John}_1) \\ & + \text{score}(\text{saw}_2, \text{movie}_4) + \text{score}(\text{saw}_2, \text{today}_5) \\ & + \text{score}(\text{movie}_4, \text{a}_3) + \dots \end{aligned}$$

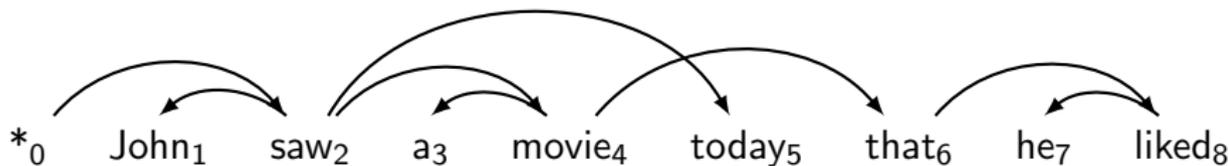
## Arc-Factored



$$\begin{aligned} f(y) = & \text{score}(\text{head} = *_0, \text{mod} = \text{saw}_2) + \text{score}(\text{saw}_2, \text{John}_1) \\ & + \text{score}(\text{saw}_2, \text{movie}_4) + \text{score}(\text{saw}_2, \text{today}_5) \\ & + \text{score}(\text{movie}_4, \text{a}_3) + \dots \end{aligned}$$

e.g.  $\text{score}(*_0, \text{saw}_2) = \log p(\text{saw}_2 | *_0)$  (generative model)

## Arc-Factored

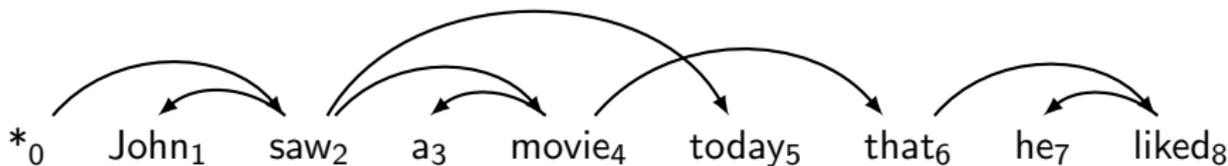


$$f(y) = \text{score}(\text{head} = *_0, \text{mod} = \text{saw}_2) + \text{score}(\text{saw}_2, \text{John}_1) \\ + \text{score}(\text{saw}_2, \text{movie}_4) + \text{score}(\text{saw}_2, \text{today}_5) \\ + \text{score}(\text{movie}_4, \text{a}_3) + \dots$$

e.g.  $\text{score}(*_0, \text{saw}_2) = \log p(\text{saw}_2 | *_0)$  (generative model)

or  $\text{score}(*_0, \text{saw}_2) = w \cdot \phi(\text{saw}_2, *_0)$  (CRF/perceptron model)

## Arc-Factored



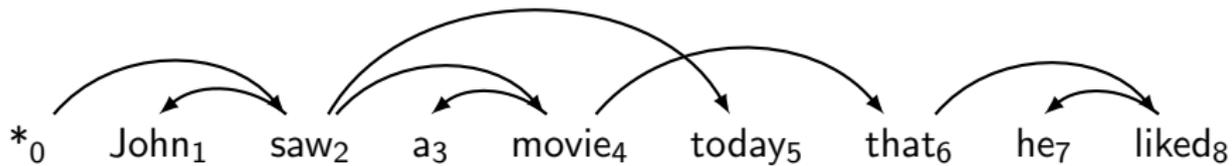
$$f(y) = \text{score}(\text{head} = *_{0}, \text{mod} = \text{saw}_{2}) + \text{score}(\text{saw}_{2}, \text{John}_{1}) \\ + \text{score}(\text{saw}_{2}, \text{movie}_{4}) + \text{score}(\text{saw}_{2}, \text{today}_{5}) \\ + \text{score}(\text{movie}_{4}, \text{a}_{3}) + \dots$$

e.g.  $\text{score}(*_{0}, \text{saw}_{2}) = \log p(\text{saw}_{2} | *_{0})$  (generative model)

or  $\text{score}(*_{0}, \text{saw}_{2}) = w \cdot \phi(\text{saw}_{2}, *_{0})$  (CRF/perceptron model)

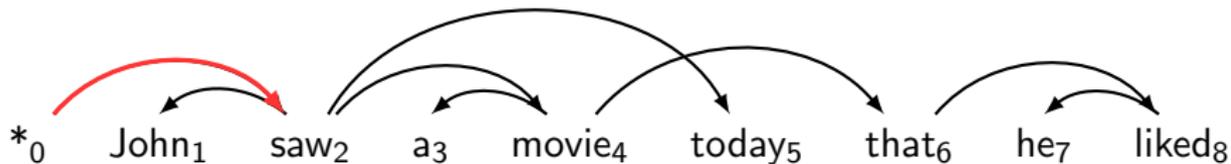
$y^* = \arg \max_y f(y) \Leftarrow$  Minimum Spanning Tree Algorithm

## Sibling Models



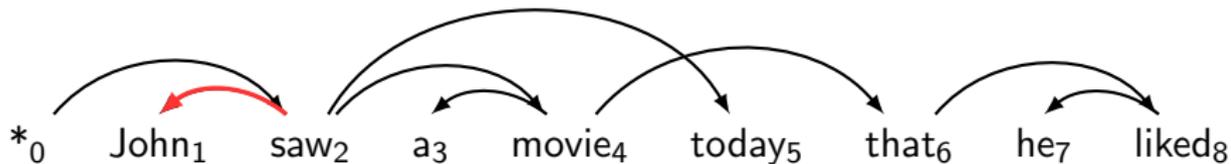
$f(y) =$

## Sibling Models



$$f(y) = \text{score}(\text{head} = *_0, \text{prev} = \text{NULL}, \text{mod} = \text{saw}_2)$$

## Sibling Models



$f(y) = \text{score}(\text{head} = *_0, \text{prev} = \text{NULL}, \text{mod} = \text{saw}_2)$

$+ \text{score}(\text{saw}_2, \text{NULL}, \text{John}_1)$

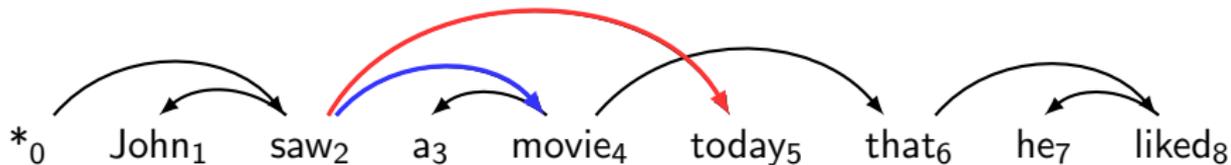
## Sibling Models



$$f(y) = \text{score}(\text{head} = *_0, \text{prev} = \text{NULL}, \text{mod} = \text{saw}_2)$$

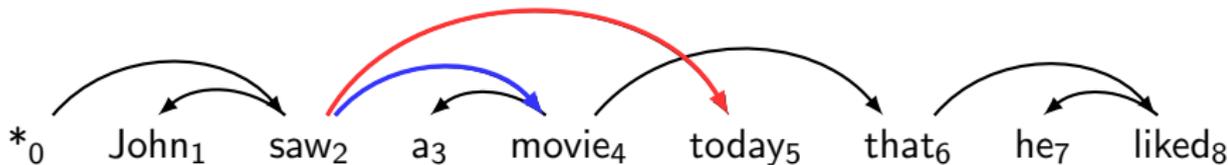
$$+ \text{score}(\text{saw}_2, \text{NULL}, \text{John}_1) + \text{score}(\text{saw}_2, \text{NULL}, \text{movie}_4)$$

## Sibling Models



$$\begin{aligned} f(y) = & \text{score}(\text{head} = *_0, \text{prev} = \text{NULL}, \text{mod} = \text{saw}_2) \\ & + \text{score}(\text{saw}_2, \text{NULL}, \text{John}_1) + \text{score}(\text{saw}_2, \text{NULL}, \text{movie}_4) \\ & + \text{score}(\text{saw}_2, \text{movie}_4, \text{today}_5) + \dots \end{aligned}$$

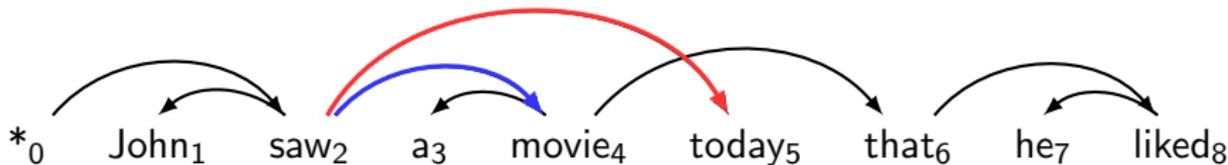
## Sibling Models



$$\begin{aligned} f(y) = & \text{score}(\text{head} = *_0, \text{prev} = \text{NULL}, \text{mod} = \text{saw}_2) \\ & + \text{score}(\text{saw}_2, \text{NULL}, \text{John}_1) + \text{score}(\text{saw}_2, \text{NULL}, \text{movie}_4) \\ & + \text{score}(\text{saw}_2, \text{movie}_4, \text{today}_5) + \dots \end{aligned}$$

e.g.  $\text{score}(\text{saw}_2, \text{movie}_4, \text{today}_5) = \log p(\text{today}_5 | \text{saw}_2, \text{movie}_4)$

## Sibling Models

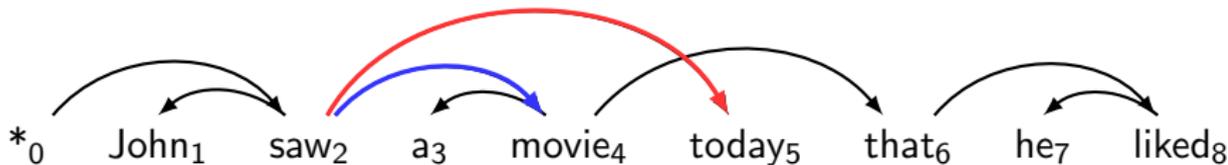


$$\begin{aligned} f(y) = & \text{score}(\text{head} = *_0, \text{prev} = \text{NULL}, \text{mod} = \text{saw}_2) \\ & + \text{score}(\text{saw}_2, \text{NULL}, \text{John}_1) + \text{score}(\text{saw}_2, \text{NULL}, \text{movie}_4) \\ & + \text{score}(\text{saw}_2, \text{movie}_4, \text{today}_5) + \dots \end{aligned}$$

e.g.  $\text{score}(\text{saw}_2, \text{movie}_4, \text{today}_5) = \log p(\text{today}_5 | \text{saw}_2, \text{movie}_4)$

or  $\text{score}(\text{saw}_2, \text{movie}_4, \text{today}_5) = w \cdot \phi(\text{saw}_2, \text{movie}_4, \text{today}_5)$

## Sibling Models



$$\begin{aligned} f(y) = & \text{score}(\text{head} = *_0, \text{prev} = \text{NULL}, \text{mod} = \text{saw}_2) \\ & + \text{score}(\text{saw}_2, \text{NULL}, \text{John}_1) + \text{score}(\text{saw}_2, \text{NULL}, \text{movie}_4) \\ & + \text{score}(\text{saw}_2, \text{movie}_4, \text{today}_5) + \dots \end{aligned}$$

e.g.  $\text{score}(\text{saw}_2, \text{movie}_4, \text{today}_5) = \log p(\text{today}_5 | \text{saw}_2, \text{movie}_4)$

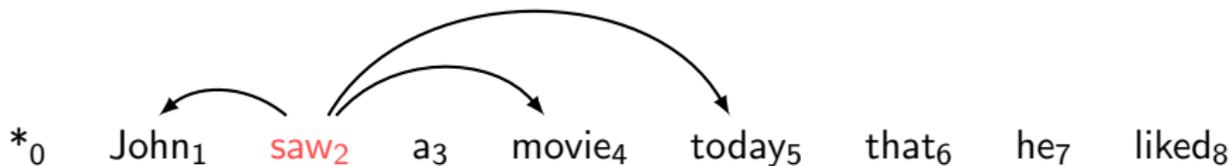
or  $\text{score}(\text{saw}_2, \text{movie}_4, \text{today}_5) = w \cdot \phi(\text{saw}_2, \text{movie}_4, \text{today}_5)$

$$y^* = \arg \max_y f(y) \Leftarrow \text{NP-Hard}$$

## Thought Experiment: Individual Decoding

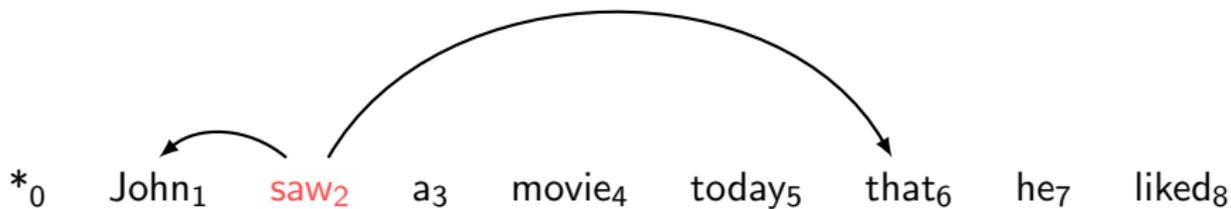
\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>

## Thought Experiment: Individual Decoding



$$\text{score}(\text{saw}_2, \text{NULL}, \text{John}_1) + \text{score}(\text{saw}_2, \text{NULL}, \text{movie}_4) \\ + \text{score}(\text{saw}_2, \text{movie}_4, \text{today}_5)$$

## Thought Experiment: Individual Decoding

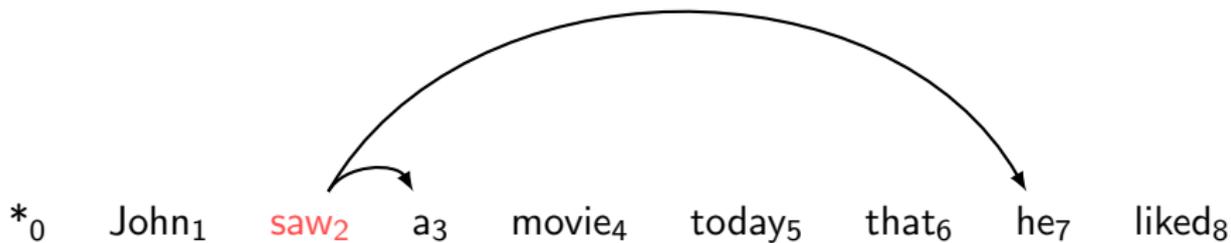


$$\text{score}(\text{saw}_2, \text{NULL}, \text{John}_1) + \text{score}(\text{saw}_2, \text{NULL}, \text{movie}_4) \\ + \text{score}(\text{saw}_2, \text{movie}_4, \text{today}_5)$$

---

$$\text{score}(\text{saw}_2, \text{NULL}, \text{John}_1) + \text{score}(\text{saw}_2, \text{NULL}, \text{that}_6)$$

## Thought Experiment: Individual Decoding



$$\text{score}(\text{saw}_2, \text{NULL}, \text{John}_1) + \text{score}(\text{saw}_2, \text{NULL}, \text{movie}_4) \\ + \text{score}(\text{saw}_2, \text{movie}_4, \text{today}_5)$$

---

$$\text{score}(\text{saw}_2, \text{NULL}, \text{John}_1) + \text{score}(\text{saw}_2, \text{NULL}, \text{that}_6)$$

---

$$\text{score}(\text{saw}_2, \text{NULL}, \text{a}_3) + \text{score}(\text{saw}_2, \text{a}_3, \text{he}_7)$$

## Thought Experiment: Individual Decoding

\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>



$2^{n-1}$   
possibilities

$score(saw_2, NULL, John_1) + score(saw_2, NULL, movie_4)$   
 $+ score(saw_2, movie_4, today_5)$

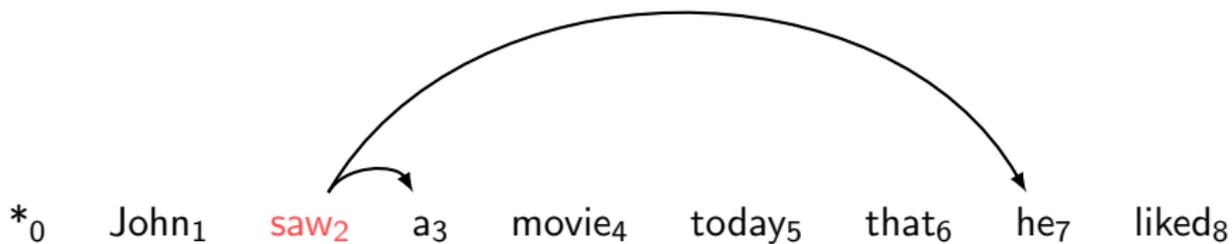
---

$score(saw_2, NULL, John_1) + score(saw_2, NULL, that_6)$

---

$score(saw_2, NULL, a_3) + score(saw_2, a_3, he_7)$

## Thought Experiment: Individual Decoding



$2^{n-1}$   
possibilities

$score(\text{saw}_2, \text{NULL}, \text{John}_1) + score(\text{saw}_2, \text{NULL}, \text{movie}_4)$   
 $+ score(\text{saw}_2, \text{movie}_4, \text{today}_5)$

---

$score(\text{saw}_2, \text{NULL}, \text{John}_1) + score(\text{saw}_2, \text{NULL}, \text{that}_6)$

---

$score(\text{saw}_2, \text{NULL}, \text{a}_3) + score(\text{saw}_2, \text{a}_3, \text{he}_7)$

Under Sibling Model, can solve for each word with [Viterbi decoding](#).

## Thought Experiment Continued

\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>



Idea: Do individual decoding for each head word using dynamic programming.

If we're lucky, we'll end up with a valid final tree.

## Thought Experiment Continued

\*<sub>0</sub>   John<sub>1</sub>   saw<sub>2</sub>   a<sub>3</sub>   movie<sub>4</sub>   today<sub>5</sub>   that<sub>6</sub>   he<sub>7</sub>   liked<sub>8</sub>



Idea: Do individual decoding for each head word using dynamic programming.

If we're lucky, we'll end up with a valid final tree.

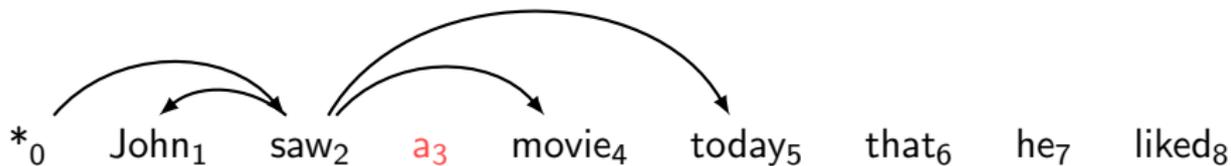
## Thought Experiment Continued



Idea: Do individual decoding for each head word using dynamic programming.

If we're lucky, we'll end up with a valid final tree.

## Thought Experiment Continued



Idea: Do individual decoding for each head word using dynamic programming.

If we're lucky, we'll end up with a valid final tree.

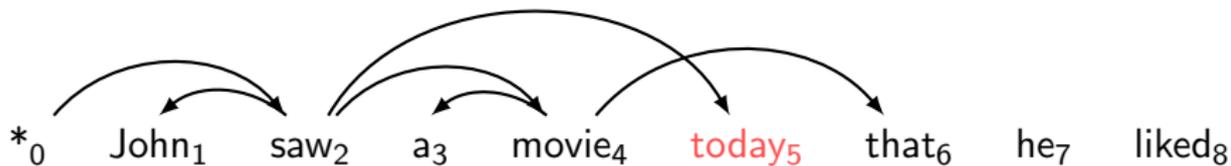
## Thought Experiment Continued



Idea: Do individual decoding for each head word using dynamic programming.

If we're lucky, we'll end up with a valid final tree.

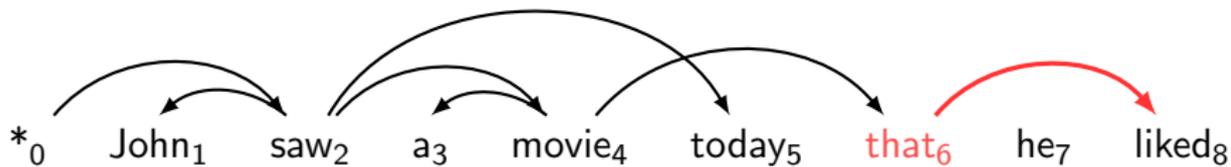
## Thought Experiment Continued



Idea: Do individual decoding for each head word using dynamic programming.

If we're lucky, we'll end up with a valid final tree.

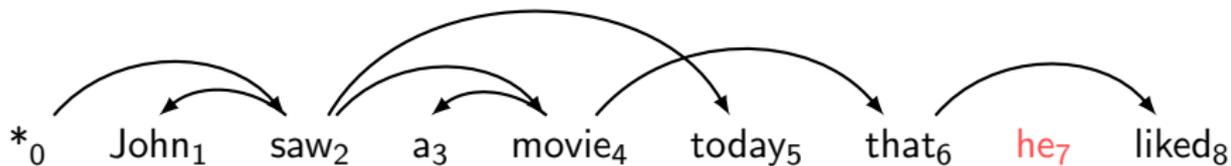
## Thought Experiment Continued



Idea: Do individual decoding for each head word using dynamic programming.

If we're lucky, we'll end up with a valid final tree.

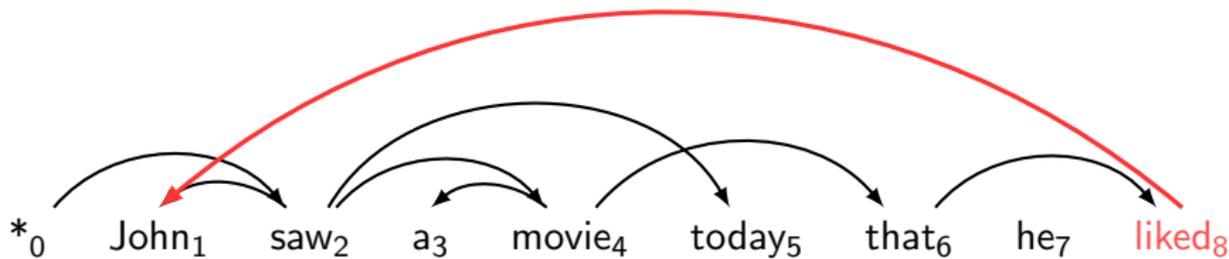
## Thought Experiment Continued



Idea: Do individual decoding for each head word using dynamic programming.

If we're lucky, we'll end up with a valid final tree.

## Thought Experiment Continued



Idea: Do individual decoding for each head word using dynamic programming.

If we're lucky, we'll end up with a valid final tree.

But we might **violate** some constraints.

# Dual Decomposition Idea

	No Constraints	Tree Constraints
Arc-Factored		Minimum Spanning Tree
Sibling Model	Individual Decoding	

# Dual Decomposition Idea

	No Constraints	Tree Constraints
Arc-Factored		Minimum Spanning Tree
Sibling Model	Individual Decoding	Dual Decomposition

# Dual Decomposition Structure

$$\text{Goal } y^* = \arg \max_{y \in \mathcal{Y}} f(y)$$

# Dual Decomposition Structure

$$\text{Goal } y^* = \arg \max_{y \in \mathcal{Y}} f(y)$$

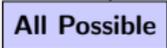
$$\text{Rewrite as } \arg \max_{z \in \mathcal{Z}, y \in \mathcal{Y}} f(z) + g(y)$$

such that  $z = y$

# Dual Decomposition Structure

$$\text{Goal } y^* = \arg \max_{y \in \mathcal{Y}} f(y)$$

Rewrite as  $\arg \max_{z \in \mathcal{Z}, y \in \mathcal{Y}} f(z) + g(y)$

All Possible

such that  $z = y$

# Dual Decomposition Structure

$$\text{Goal } y^* = \arg \max_{y \in \mathcal{Y}} f(y)$$

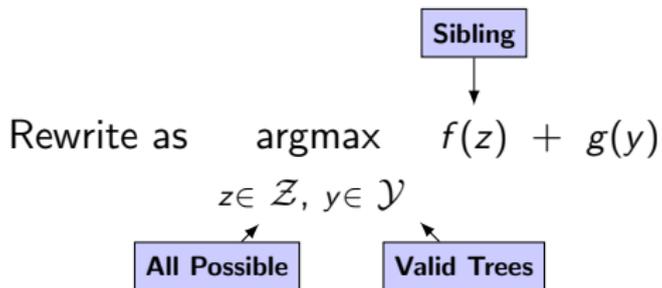
Rewrite as  $\arg \max_{z \in \mathcal{Z}, y \in \mathcal{Y}} f(z) + g(y)$



such that  $z = y$

# Dual Decomposition Structure

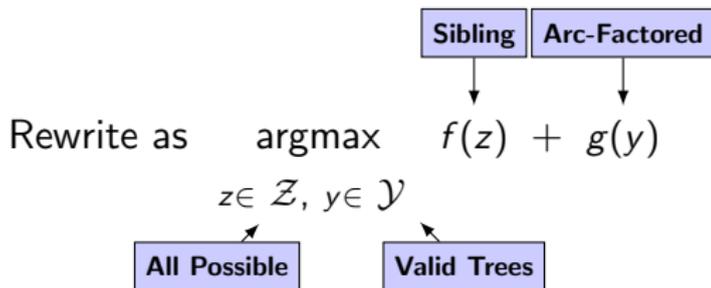
$$\text{Goal } y^* = \arg \max_{y \in \mathcal{Y}} f(y)$$



such that  $z = y$

# Dual Decomposition Structure

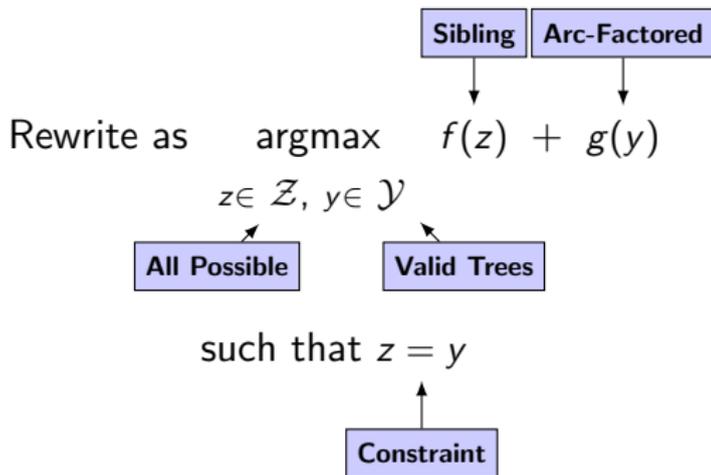
$$\text{Goal } y^* = \arg \max_{y \in \mathcal{Y}} f(y)$$



such that  $z = y$

# Dual Decomposition Structure

$$\text{Goal } y^* = \arg \max_{y \in \mathcal{Y}} f(y)$$



## Algorithm Sketch

Set penalty weights equal to 0 for all edges.

**For**  $k = 1$  **to**  $K$

## Algorithm Sketch

Set penalty weights equal to 0 for all edges.

**For**  $k = 1$  **to**  $K$

$z^{(k)} \leftarrow$  Decode  $(f(z) + \text{penalty})$  by Individual Decoding

## Algorithm Sketch

Set penalty weights equal to 0 for all edges.

**For**  $k = 1$  **to**  $K$

$z^{(k)} \leftarrow$  Decode  $(f(z) + \text{penalty})$  by Individual Decoding

$y^{(k)} \leftarrow$  Decode  $(g(y) - \text{penalty})$  by Minimum Spanning Tree

## Algorithm Sketch

Set penalty weights equal to 0 for all edges.

**For**  $k = 1$  **to**  $K$

$z^{(k)} \leftarrow$  Decode  $(f(z) + \text{penalty})$  by Individual Decoding

$y^{(k)} \leftarrow$  Decode  $(g(y) - \text{penalty})$  by Minimum Spanning Tree

**If**  $y^{(k)}(i,j) = z^{(k)}(i,j)$  for all  $i,j$  **Return**  $(y^{(k)}, z^{(k)})$

## Algorithm Sketch

Set penalty weights equal to 0 for all edges.

**For**  $k = 1$  **to**  $K$

$z^{(k)} \leftarrow$  Decode  $(f(z) + \text{penalty})$  by Individual Decoding

$y^{(k)} \leftarrow$  Decode  $(g(y) - \text{penalty})$  by Minimum Spanning Tree

**If**  $y^{(k)}(i,j) = z^{(k)}(i,j)$  for all  $i,j$  **Return**  $(y^{(k)}, z^{(k)})$

**Else** Update penalty weights based on  $y^{(k)}(i,j) - z^{(k)}(i,j)$

## Individual Decoding

## Penalties

$u(i, j) = 0$  for all  $i, j$

\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>

$$z^* = \arg \max_{z \in \mathcal{Z}} (f(z) + \sum_{i,j} u(i,j)z(i,j))$$

## Minimum Spanning Tree

\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>

$$y^* = \arg \max_{y \in \mathcal{Y}} (g(y) - \sum_{i,j} u(i,j)y(i,j))$$

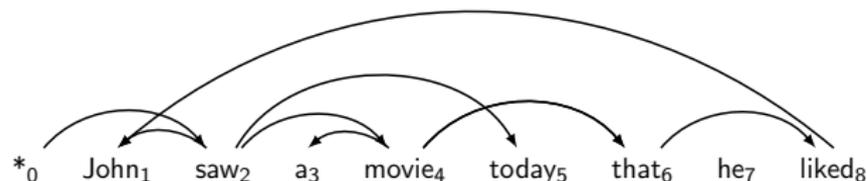
## Key

$f(z)$	$\Leftarrow$	Sibling Model	$g(y)$	$\Leftarrow$	Arc-Factored Model
$\mathcal{Z}$	$\Leftarrow$	No Constraints	$\mathcal{Y}$	$\Leftarrow$	Tree Constraints
$y(i,j) = 1$	if	$y$ contains dependency $i,j$			

## Individual Decoding

## Penalties

$$u(i,j) = 0 \text{ for all } i,j$$



$$z^* = \arg \max_{z \in \mathcal{Z}} (f(z) + \sum_{i,j} u(i,j)z(i,j))$$

## Minimum Spanning Tree

\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>

$$y^* = \arg \max_{y \in \mathcal{Y}} (g(y) - \sum_{i,j} u(i,j)y(i,j))$$

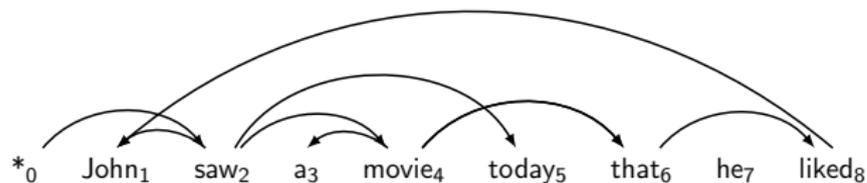
## Key

$f(z)$	$\Leftarrow$	Sibling Model	$g(y)$	$\Leftarrow$	Arc-Factored Model
$\mathcal{Z}$	$\Leftarrow$	No Constraints	$\mathcal{Y}$	$\Leftarrow$	Tree Constraints
$y(i,j) = 1$	if	$y$ contains dependency $i,j$			

## Individual Decoding

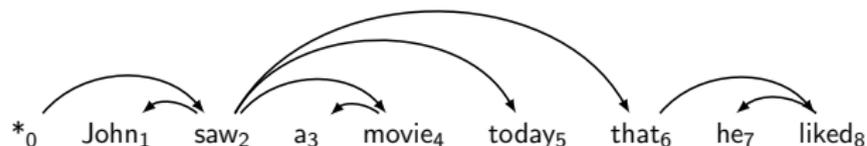
## Penalties

$$u(i,j) = 0 \text{ for all } i,j$$



$$z^* = \arg \max_{z \in \mathcal{Z}} (f(z) + \sum_{i,j} u(i,j)z(i,j))$$

## Minimum Spanning Tree



$$y^* = \arg \max_{y \in \mathcal{Y}} (g(y) - \sum_{i,j} u(i,j)y(i,j))$$

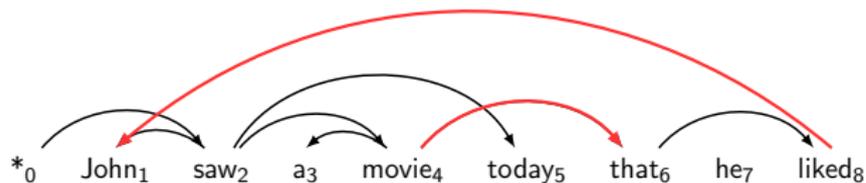
## Key

$f(z)$	$\Leftarrow$	Sibling Model	$g(y)$	$\Leftarrow$	Arc-Factored Model
$\mathcal{Z}$	$\Leftarrow$	No Constraints	$\mathcal{Y}$	$\Leftarrow$	Tree Constraints
$y(i,j) = 1$	if	$y$ contains dependency $i,j$			

## Individual Decoding

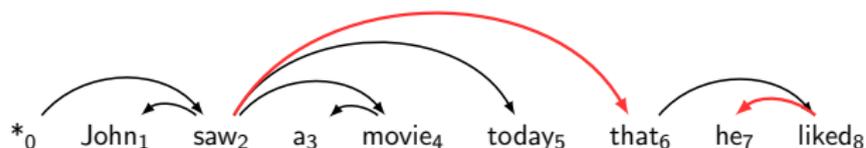
## Penalties

$$u(i,j) = 0 \text{ for all } i,j$$



$$z^* = \arg \max_{z \in \mathcal{Z}} (f(z) + \sum_{i,j} u(i,j)z(i,j))$$

## Minimum Spanning Tree

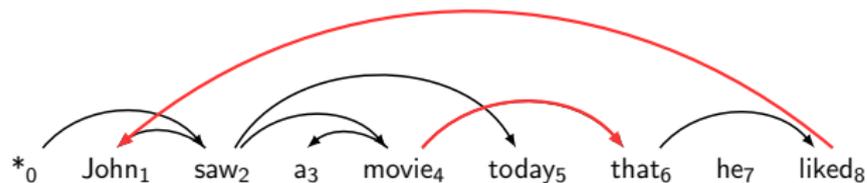


$$y^* = \arg \max_{y \in \mathcal{Y}} (g(y) - \sum_{i,j} u(i,j)y(i,j))$$

## Key

$f(z)$	$\Leftarrow$ Sibling Model	$g(y)$	$\Leftarrow$ Arc-Factored Model
$\mathcal{Z}$	$\Leftarrow$ No Constraints	$\mathcal{Y}$	$\Leftarrow$ Tree Constraints
$y(i,j) = 1$	if $y$ contains dependency $i,j$		

## Individual Decoding



$$z^* = \arg \max_{z \in \mathcal{Z}} (f(z) + \sum_{i,j} u(i,j)z(i,j))$$

## Penalties

$u(i,j) = 0$  for all  $i,j$

Iteration 1

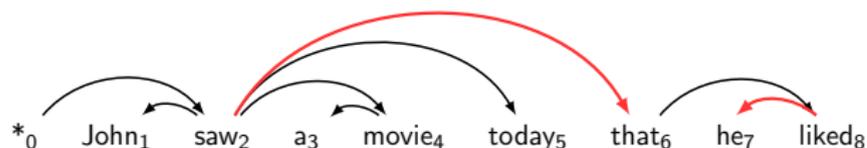
$u(8,1) = -1$

$u(4,6) = -1$

$u(2,6) = 1$

$u(8,7) = 1$

## Minimum Spanning Tree

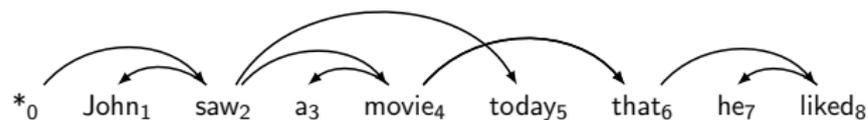


$$y^* = \arg \max_{y \in \mathcal{Y}} (g(y) - \sum_{i,j} u(i,j)y(i,j))$$

## Key

$f(z)$	$\Leftarrow$ Sibling Model	$g(y)$	$\Leftarrow$ Arc-Factored Model
$\mathcal{Z}$	$\Leftarrow$ No Constraints	$\mathcal{Y}$	$\Leftarrow$ Tree Constraints
$y(i,j) = 1$	if $y$ contains dependency $i,j$		

## Individual Decoding



$$z^* = \arg \max_{z \in \mathcal{Z}} (f(z) + \sum_{i,j} u(i,j)z(i,j))$$

## Penalties

$u(i,j) = 0$  for all  $i,j$

Iteration 1

$u(8,1) = -1$

$u(4,6) = -1$

$u(2,6) = 1$

$u(8,7) = 1$

## Minimum Spanning Tree

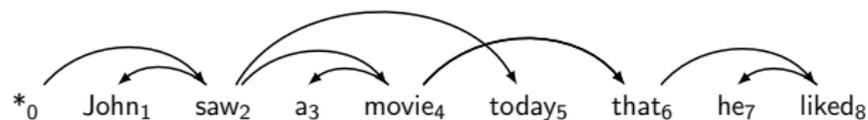
\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>

$$y^* = \arg \max_{y \in \mathcal{Y}} (g(y) - \sum_{i,j} u(i,j)y(i,j))$$

## Key

$f(z)$	$\Leftarrow$	Sibling Model	$g(y)$	$\Leftarrow$	Arc-Factored Model
$\mathcal{Z}$	$\Leftarrow$	No Constraints	$\mathcal{Y}$	$\Leftarrow$	Tree Constraints
$y(i,j) = 1$	if	$y$ contains dependency $i,j$			

## Individual Decoding



$$z^* = \arg \max_{z \in \mathcal{Z}} (f(z) + \sum_{i,j} u(i,j)z(i,j))$$

## Penalties

$u(i,j) = 0$  for all  $i,j$

Iteration 1

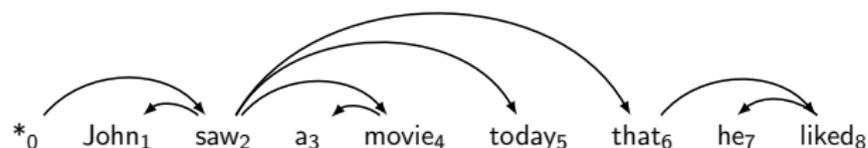
$u(8,1) = -1$

$u(4,6) = -1$

$u(2,6) = 1$

$u(8,7) = 1$

## Minimum Spanning Tree

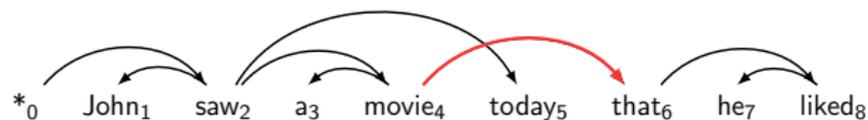


$$y^* = \arg \max_{y \in \mathcal{Y}} (g(y) - \sum_{i,j} u(i,j)y(i,j))$$

## Key

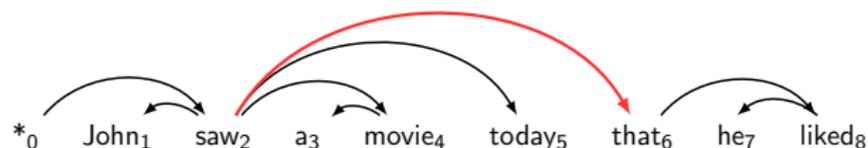
$f(z)$	$\Leftarrow$ Sibling Model	$g(y)$	$\Leftarrow$ Arc-Factored Model
$\mathcal{Z}$	$\Leftarrow$ No Constraints	$\mathcal{Y}$	$\Leftarrow$ Tree Constraints
$y(i,j) = 1$	if $y$ contains dependency $i,j$		

## Individual Decoding



$$z^* = \arg \max_{z \in \mathcal{Z}} (f(z) + \sum_{i,j} u(i,j)z(i,j))$$

## Minimum Spanning Tree



$$y^* = \arg \max_{y \in \mathcal{Y}} (g(y) - \sum_{i,j} u(i,j)y(i,j))$$

## Key

$f(z)$	$\Leftarrow$ Sibling Model	$g(y)$	$\Leftarrow$ Arc-Factored Model
$\mathcal{Z}$	$\Leftarrow$ No Constraints	$\mathcal{Y}$	$\Leftarrow$ Tree Constraints
$y(i,j) = 1$	if $y$ contains dependency $i,j$		

## Penalties

$u(i,j) = 0$  for all  $i,j$

Iteration 1

$u(8,1) \quad -1$

$u(4,6) \quad -1$

$u(2,6) \quad 1$

$u(8,7) \quad 1$

Iteration 2

$u(8,1) \quad -1$

$u(4,6) \quad -2$

$u(2,6) \quad 2$

$u(8,7) \quad 1$

## Individual Decoding

\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>

$$z^* = \arg \max_{z \in \mathcal{Z}} (f(z) + \sum_{i,j} u(i,j)z(i,j))$$

## Minimum Spanning Tree

\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>

$$y^* = \arg \max_{y \in \mathcal{Y}} (g(y) - \sum_{i,j} u(i,j)y(i,j))$$

## Key

$f(z)$	$\Leftarrow$	Sibling Model	$g(y)$	$\Leftarrow$	Arc-Factored Model
$\mathcal{Z}$	$\Leftarrow$	No Constraints	$\mathcal{Y}$	$\Leftarrow$	Tree Constraints
$y(i,j) = 1$	if	$y$ contains dependency $i,j$			

## Penalties

$u(i,j) = 0$  for all  $i,j$

Iteration 1

---

$u(8,1)$	-1
----------	----

$u(4,6)$	-1
----------	----

$u(2,6)$	1
----------	---

$u(8,7)$	1
----------	---

Iteration 2

---

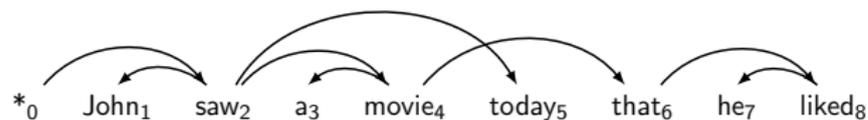
$u(8,1)$	-1
----------	----

$u(4,6)$	-2
----------	----

$u(2,6)$	2
----------	---

$u(8,7)$	1
----------	---

## Individual Decoding



$$z^* = \arg \max_{z \in \mathcal{Z}} (f(z) + \sum_{i,j} u(i,j)z(i,j))$$

## Minimum Spanning Tree

\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>

$$y^* = \arg \max_{y \in \mathcal{Y}} (g(y) - \sum_{i,j} u(i,j)y(i,j))$$

## Key

$f(z)$	$\Leftarrow$	Sibling Model	$g(y)$	$\Leftarrow$	Arc-Factored Model
$\mathcal{Z}$	$\Leftarrow$	No Constraints	$\mathcal{Y}$	$\Leftarrow$	Tree Constraints
$y(i,j) = 1$	if	$y$ contains dependency $i,j$			

## Penalties

$u(i,j) = 0$  for all  $i,j$

Iteration 1

$u(8,1) \quad -1$

$u(4,6) \quad -1$

$u(2,6) \quad 1$

$u(8,7) \quad 1$

Iteration 2

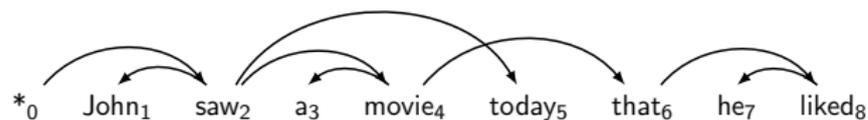
$u(8,1) \quad -1$

$u(4,6) \quad -2$

$u(2,6) \quad 2$

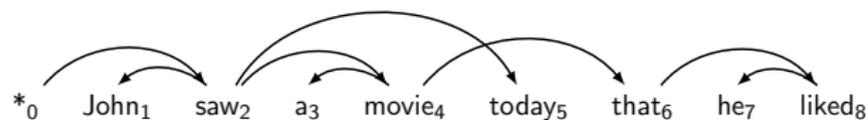
$u(8,7) \quad 1$

## Individual Decoding



$$z^* = \arg \max_{z \in \mathcal{Z}} (f(z) + \sum_{i,j} u(i,j)z(i,j))$$

## Minimum Spanning Tree



$$y^* = \arg \max_{y \in \mathcal{Y}} (g(y) - \sum_{i,j} u(i,j)y(i,j))$$

## Key

$f(z)$	$\Leftarrow$	Sibling Model	$g(y)$	$\Leftarrow$	Arc-Factored Model
$\mathcal{Z}$	$\Leftarrow$	No Constraints	$\mathcal{Y}$	$\Leftarrow$	Tree Constraints
$y(i,j) = 1$	if	$y$ contains dependency $i,j$			

## Penalties

$u(i,j) = 0$  for all  $i,j$

Iteration 1

$u(8,1) = -1$

$u(4,6) = -1$

$u(2,6) = 1$

$u(8,7) = 1$

Iteration 2

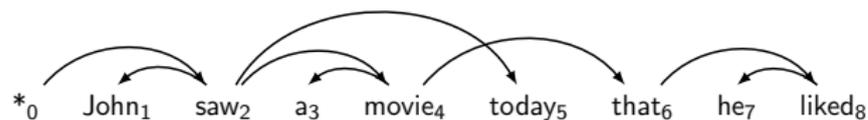
$u(8,1) = -1$

$u(4,6) = -2$

$u(2,6) = 2$

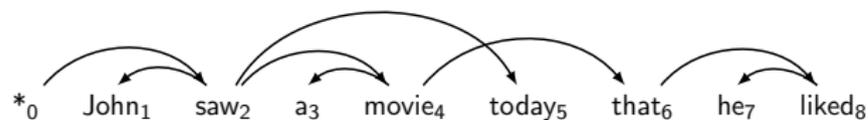
$u(8,7) = 1$

## Individual Decoding



$$z^* = \arg \max_{z \in \mathcal{Z}} (f(z) + \sum_{i,j} u(i,j)z(i,j))$$

## Minimum Spanning Tree



$$y^* = \arg \max_{y \in \mathcal{Y}} (g(y) - \sum_{i,j} u(i,j)y(i,j))$$

## Key

$f(z)$	$\Leftarrow$ Sibling Model	$g(y)$	$\Leftarrow$ Arc-Factored Model
$\mathcal{Z}$	$\Leftarrow$ No Constraints	$\mathcal{Y}$	$\Leftarrow$ Tree Constraints
$y(i,j) = 1$	if $y$ contains dependency $i,j$		

## Penalties

$u(i,j) = 0$  for all  $i,j$

Iteration 1

$u(8,1) = -1$

$u(4,6) = -1$

$u(2,6) = 1$

$u(8,7) = 1$

Iteration 2

$u(8,1) = -1$

$u(4,6) = -2$

$u(2,6) = 2$

$u(8,7) = 1$

**Converged**

$$y^* = \arg \max_{y \in \mathcal{Y}} f(y) + g(y)$$

# Guarantees

## Theorem

If at any iteration  $y^{(k)} = z^{(k)}$ , then  $(y^{(k)}, z^{(k)})$  is the global optimum.

In experiments, we find the global optimum on 98% of examples.

# Guarantees

## Theorem

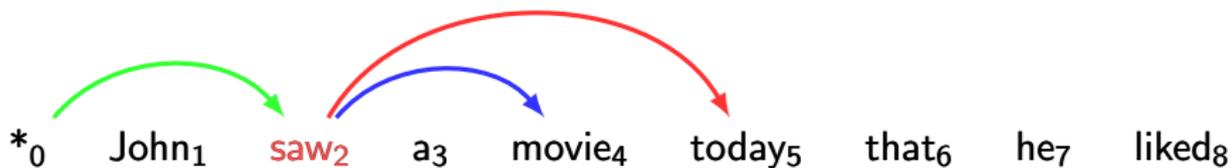
If at any iteration  $y^{(k)} = z^{(k)}$ , then  $(y^{(k)}, z^{(k)})$  is the global optimum.

In experiments, we find the global optimum on 98% of examples.

If we do not converge to a match, we can still return an approximate solution (more in the paper).

# Extensions

## ▶ Grandparent Models



$$f(y) = \dots + \text{score}(gp = *_0, head = \text{saw}_2, prev = \text{movie}_4, mod = \text{today}_5)$$

## ▶ Head Automata (Eisner, 2000)

Generalization of Sibling models

Allow arbitrary automata as local scoring function.

Roadmap

Algorithm

Experiments

Derivation

# Experiments

Properties:

- ▶ Exactness
- ▶ Parsing Speed
- ▶ Parsing Accuracy
- ▶ Comparison to Individual Decoding
- ▶ Comparison to LP/ILP

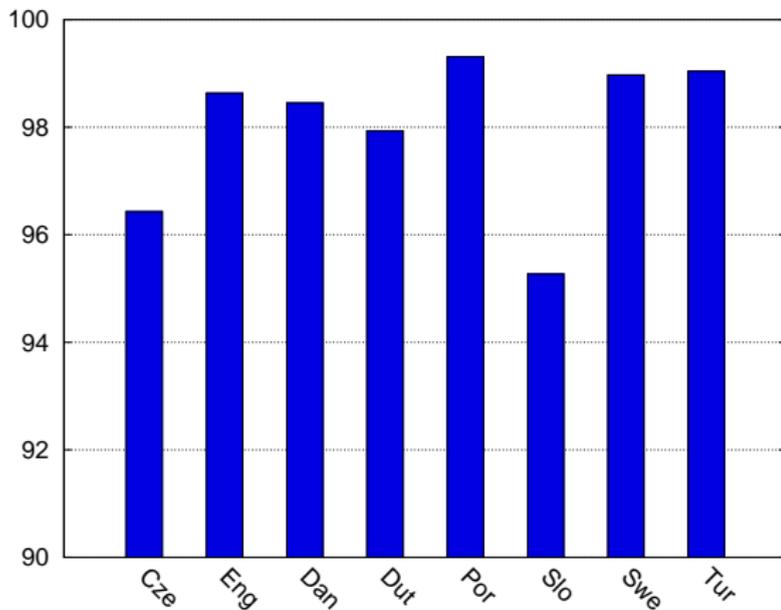
Training:

- ▶ Averaged Perceptron (more details in paper)

Experiments on:

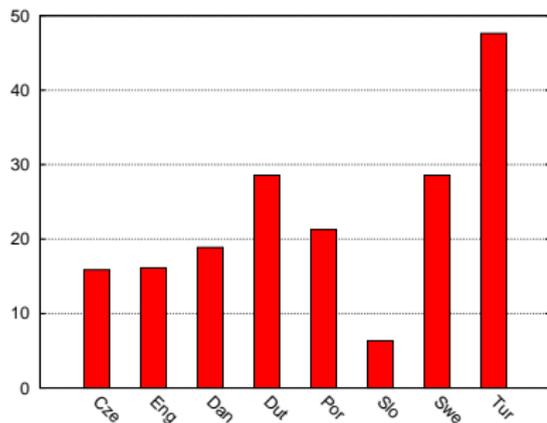
- ▶ CoNLL Datasets
- ▶ English Penn Treebank
- ▶ Czech Dependency Treebank

## How often do we exactly solve the problem?

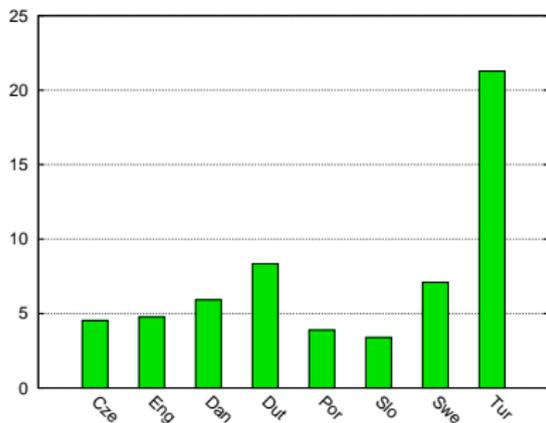


- ▶ Percentage of examples where the dual decomposition finds an exact solution.

# Parsing Speed



Sibling model



Grandparent model

- ▶ Number of sentences parsed per second
- ▶ Comparable to dynamic programming for projective parsing

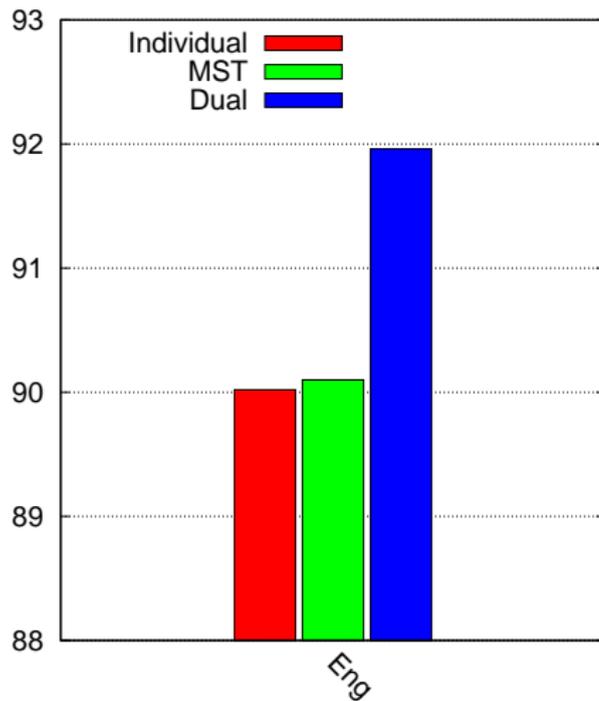
# Accuracy

	Arc-Factored	Prev Best	Grandparent
Dan	89.7	91.5	<b>91.8</b>
Dut	82.3	85.6	<b>85.8</b>
Por	90.7	92.1	<b>93.0</b>
Slo	82.4	85.6	<b>86.2</b>
Swe	88.9	90.6	<b>91.4</b>
Tur	75.7	76.4	<b>77.6</b>
Eng	90.1	—	<b>92.5</b>
Cze	84.4	—	<b>87.3</b>

Prev Best - Best reported results for CoNLL-X data set, includes

- ▶ Approximate search (McDonald and Pereira, 2006)
- ▶ Loop belief propagation (Smith and Eisner, 2008)
- ▶ (Integer) Linear Programming (Martins et al., 2009)

## Comparison to Subproblems



F<sub>1</sub> for dependency accuracy

## Comparison to LP/ILP

Martins et al.(2009): Proposes two representations of non-projective dependency parsing as a linear programming relaxation as well as an exact ILP.

- ▶ LP (1)
- ▶ LP (2)
- ▶ ILP

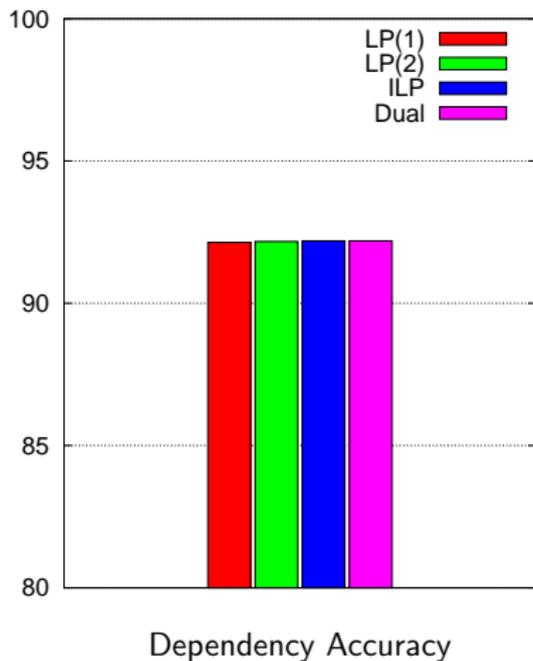
Use an LP/ILP Solver for decoding

We compare:

- ▶ Accuracy
- ▶ Exactness
- ▶ Speed

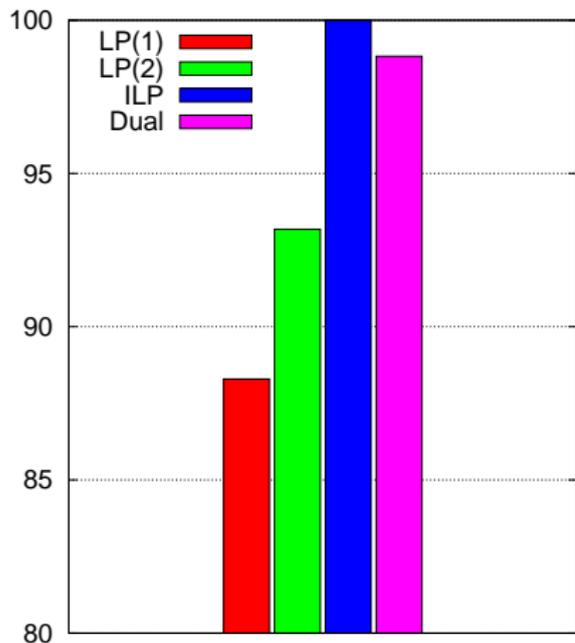
Both LP and dual decomposition methods use the same model, features, and weights  $w$ .

## Comparison to LP/ILP: Accuracy

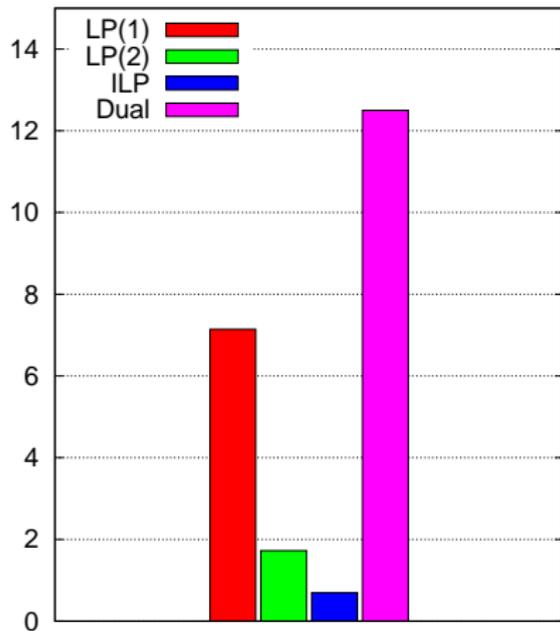


- ▶ All decoding methods have comparable accuracy

# Comparison to LP/ILP: Exactness and Speed



Percentage with exact solution



Sentences per second

Roadmap

Algorithm

Experiments

Derivation

## Deriving the Algorithm

Goal:

$$y^* = \arg \max_{y \in \mathcal{Y}} f(y)$$

Rewrite:

$$\arg \max_{z \in \mathcal{Z}, y \in \mathcal{Y}} f(z) + g(y)$$

s.t.  $z(i,j) = y(i,j)$  for all  $i,j$

$$\text{Lagrangian: } L(u, y, z) = f(z) + g(y) + \sum_{i,j} u(i,j) (z(i,j) - y(i,j))$$

## Deriving the Algorithm

Goal:

$$y^* = \arg \max_{y \in \mathcal{Y}} f(y)$$

Rewrite:

$$\arg \max_{z \in \mathcal{Z}, y \in \mathcal{Y}} f(z) + g(y)$$

$$\text{s.t. } z(i, j) = y(i, j) \text{ for all } i, j$$

$$\text{Lagrangian: } L(u, y, z) = f(z) + g(y) + \sum_{i,j} u(i, j) (z(i, j) - y(i, j))$$

The **dual problem** is to find  $\min_u L(u)$  where

$$\begin{aligned} L(u) = \max_{y \in \mathcal{Y}, z \in \mathcal{Z}} L(u, y, z) &= \max_{z \in \mathcal{Z}} \left( f(z) + \sum_{i,j} u(i, j) z(i, j) \right) \\ &+ \max_{y \in \mathcal{Y}} \left( g(y) - \sum_{i,j} u(i, j) y(i, j) \right) \end{aligned}$$

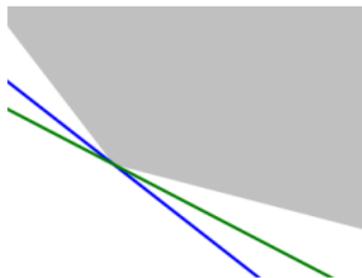
Dual is an upper bound:  $L(u) \geq f(z^*) + g(y^*)$  for any  $u$

## A Subgradient Algorithm for Minimizing $L(u)$

$$L(u) = \max_{z \in \mathcal{Z}} \left( f(z) + \sum_{i,j} u(i,j)z(i,j) \right) + \max_{y \in \mathcal{Y}} \left( g(y) - \sum_{i,j} u(i,j)y(i,j) \right)$$

$L(u)$  is convex, but not differentiable. A *subgradient* of  $L(u)$  at  $u$  is a vector  $g_u$  such that for all  $v$ ,

$$L(v) \geq L(u) + g_u \cdot (v - u)$$



Subgradient methods use updates  $u' = u - \alpha g_u$

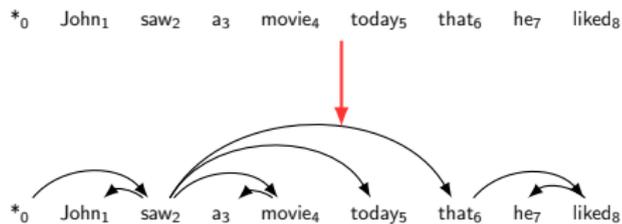
In fact, for our  $L(u)$ ,  $g_u(i,j) = z^*(i,j) - y^*(i,j)$

## Related Work

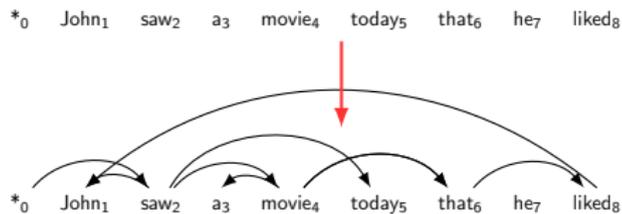
- ▶ Methods that use general purpose linear programming or integer linear programming solvers (Martins et al. 2009; Riedel and Clarke 2006; Roth and Yih 2005)
- ▶ Dual decomposition/Lagrangian relaxation in combinatorial optimization (Dantzig and Wolfe, 1960; Held and Karp, 1970; Fisher 1981)
- ▶ Dual decomposition for inference in MRFs (Komodakis et al., 2007; Wainwright et al., 2005)
- ▶ Methods that incorporate combinatorial solvers within loopy belief propagation (Duchi et al. 2007; Smith and Eisner 2008)

# Summary

$$y^* = \arg \max_y f(y) \Leftarrow \text{NP-Hard}$$



Arc-Factored Model

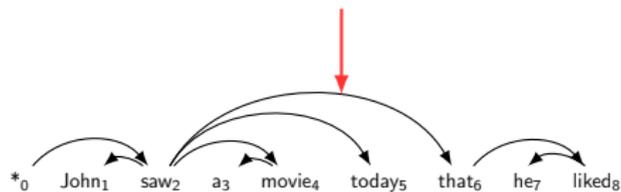


Sibling Model

# Summary

$$y^* = \arg \max_y f(y) \Leftarrow \text{NP-Hard}$$

\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>

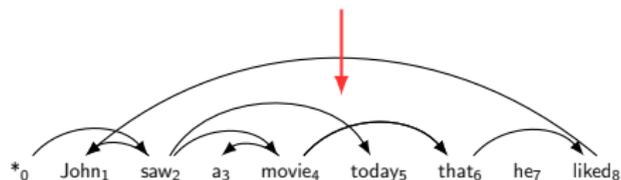


Arc-Factored Model

+

Dual Decomposition

\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>



Sibling Model

## Other Applications

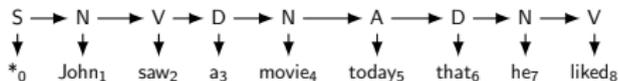
- ▶ Dual decomposition can be applied to other decoding problems.
- ▶ Rush et al. (2010) focuses on integrated dynamic programming algorithms.
  - ▶ Integrated Parsing and Tagging
  - ▶ Integrated Constituency and Dependency Parsing



# Parsing and Tagging

$$y^* = \arg \max_y f(y) \Leftarrow \text{Slow}$$

\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>

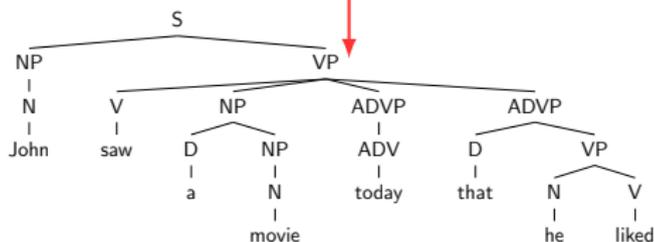


HMM Model



Dual Decomposition

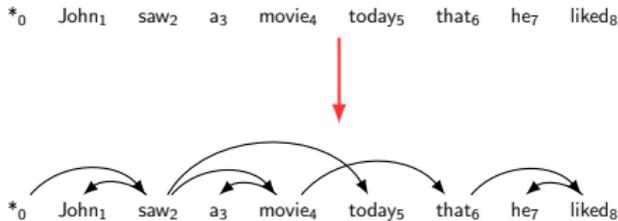
\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>



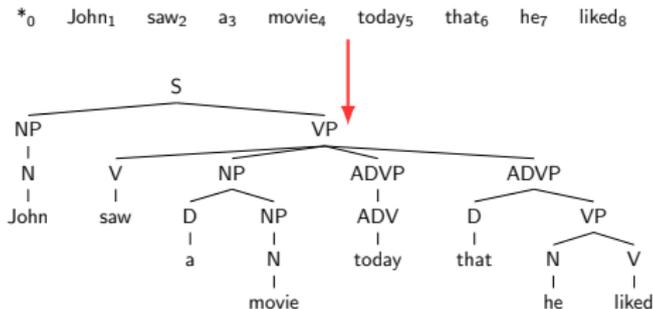
CFG Model

# Dependency and Constituency

$$y^* = \arg \max_y f(y) \Leftarrow \text{Slow}$$



Dependency Model



Lexicalized CFG

# Dependency and Constituency

$$y^* = \arg \max_y f(y) \Leftarrow \text{Slow}$$

\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>

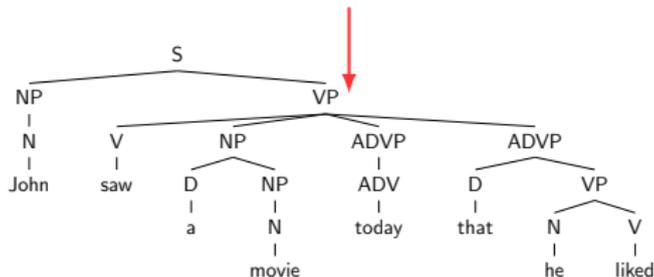


Dependency Model



Dual Decomposition

\*<sub>0</sub> John<sub>1</sub> saw<sub>2</sub> a<sub>3</sub> movie<sub>4</sub> today<sub>5</sub> that<sub>6</sub> he<sub>7</sub> liked<sub>8</sub>



Lexicalized CFG

# Future Directions

There is much more to explore around dual decomposition in NLP.

- ▶ Known Techniques

- ▶ Generalization to more than two models
- ▶ K-best decoding
- ▶ Approximate subgradient
- ▶ Heuristic for branch-and-bound type search

- ▶ Possible NLP Applications

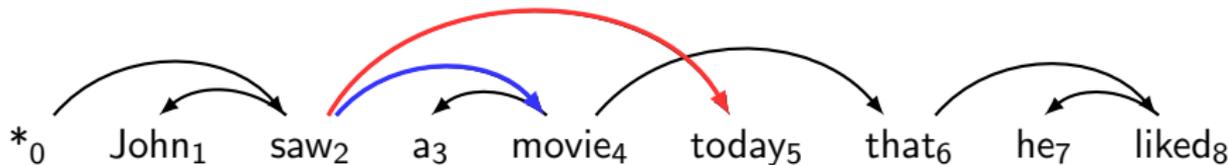
- ▶ Machine Translation
- ▶ Speech Recognition
- ▶ “Loopy” Sequence Models

- ▶ Open Questions

- ▶ Can we speed up subalgorithms when running repeatedly?
- ▶ What are the trade-offs of different decompositions?
- ▶ Are there better methods for optimizing the dual?

## Appendix

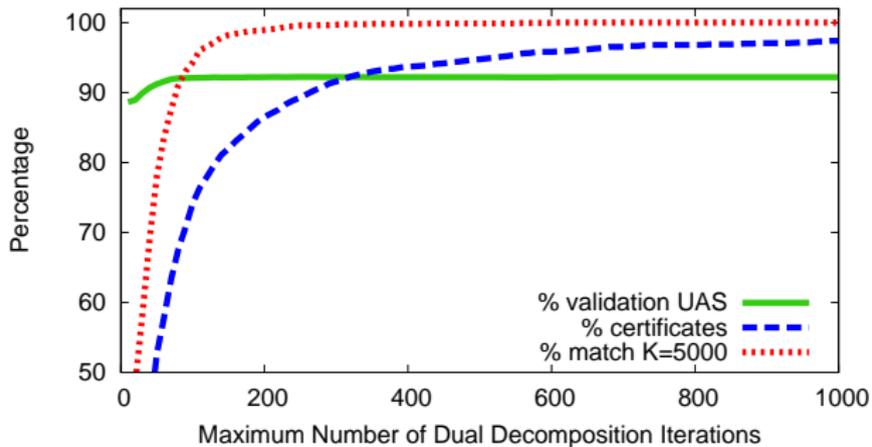
## Training the Model



$$f(y) = \dots + \text{score}(\text{saw}_2, \text{movie}_4, \text{today}_5) + \dots$$

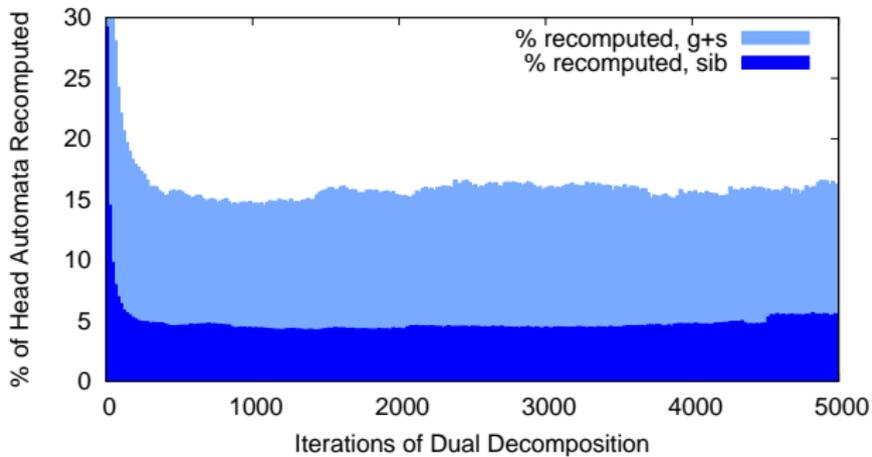
- ▶  $\text{score}(\text{saw}_2, \text{movie}_4, \text{today}_5) = w \cdot \phi(\text{saw}_2, \text{movie}_4, \text{today}_5)$
- ▶ Weight vector  $w$  trained using **Averaged perceptron**.
- ▶ (More details in the paper.)

# Early Stopping



Early Stopping

# Caching



Caching speed