

# *K*-best Parsing: Algorithms and Applications



Liang Huang 黃亮  
University of Pennsylvania

*joint work with David Chiang 蔣偉*  
(USC Information Sciences Institute)

---

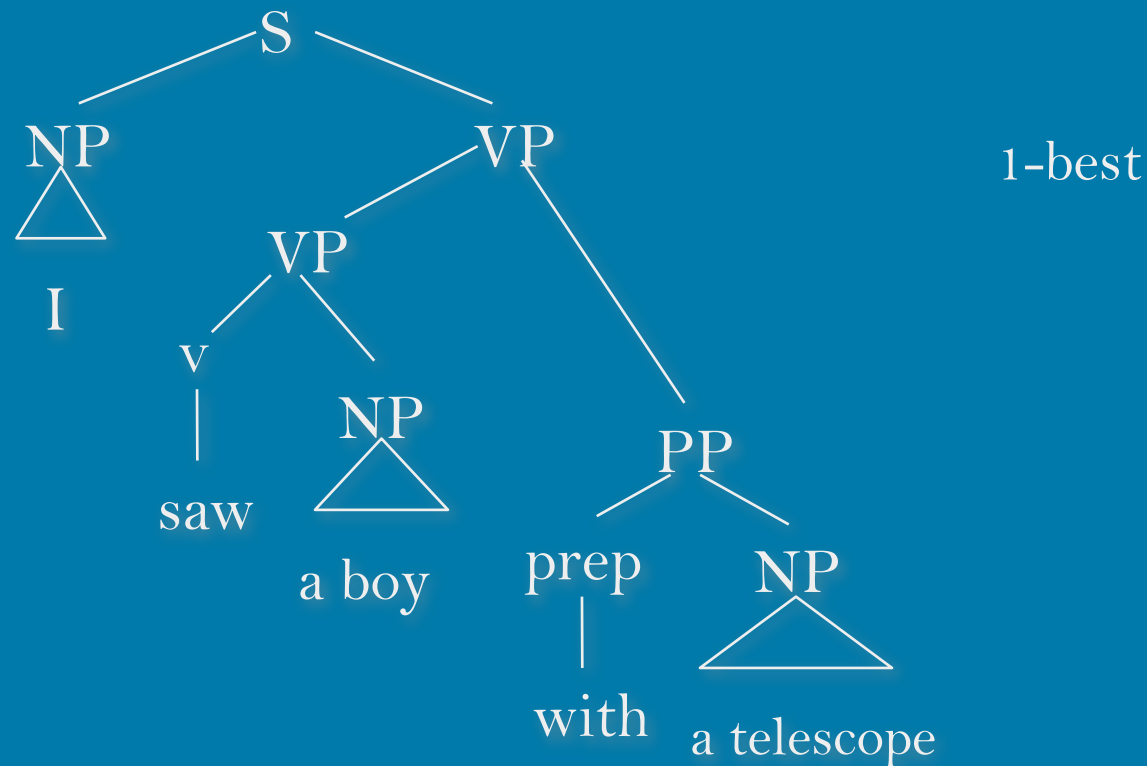
Hong Kong University of Science & Technology  
July 30, 2007

# *k*-best Parsing

# $k$ -best Parsing

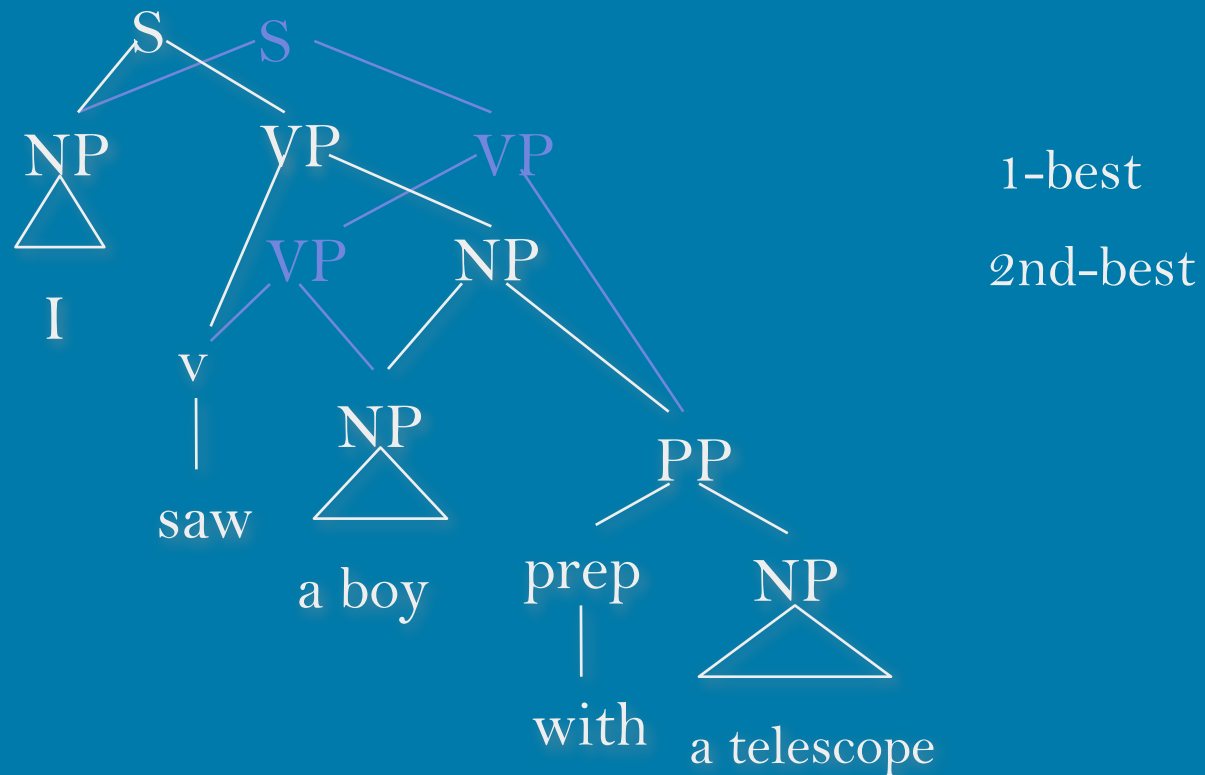
I saw a boy with a telescope.

# $k$ -best Parsing



I saw a boy with a telescope.

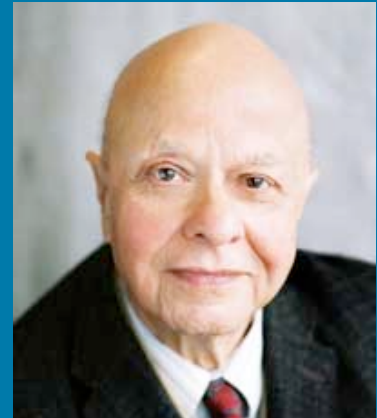
# $k$ -best Parsing



I saw a boy with a telescope.

# Not a trivial task...

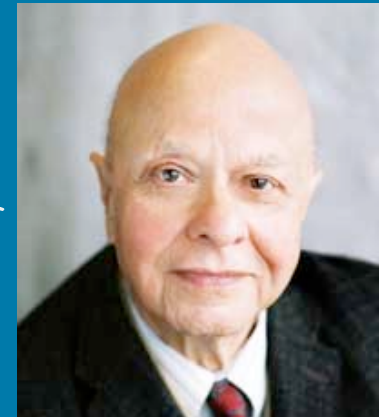
# Not a trivial task...



Aravind Joshi

# Not a trivial task...

I saw her duck.

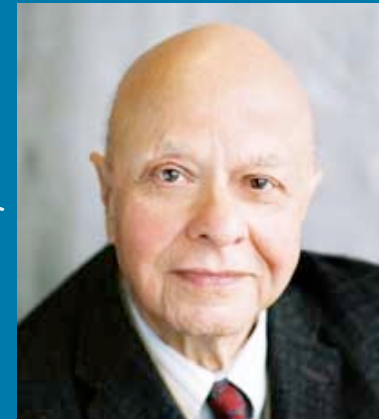


Aravind Joshi



# Not a trivial task...

I saw her duck.

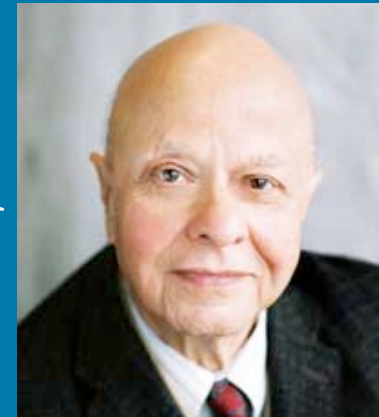


Aravind Joshi

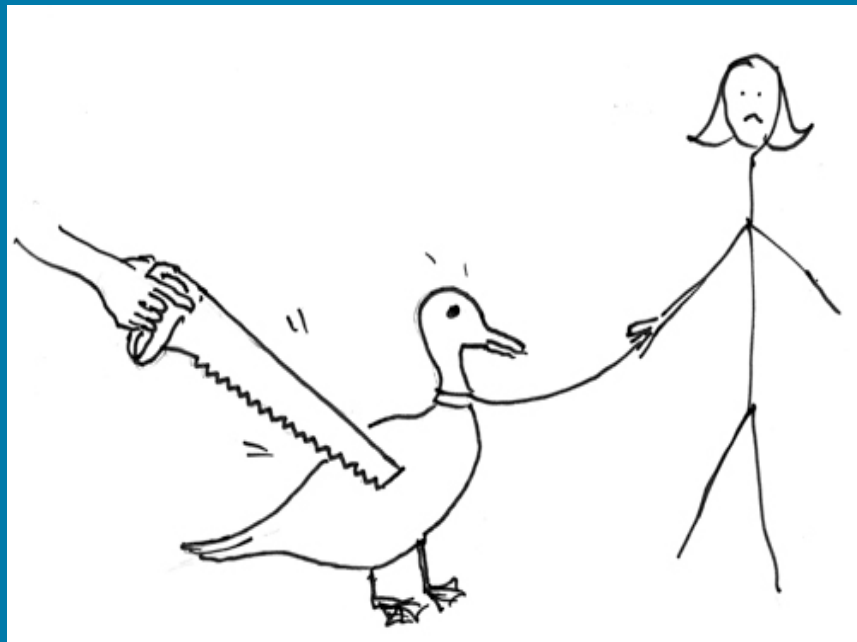


# Not a trivial task...

I saw her duck.

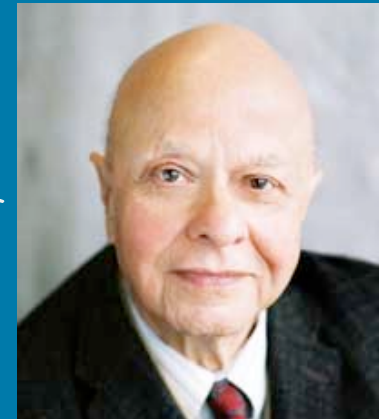


Aravind Joshi



# Not a trivial task...

I eat sushi with tuna.



Aravind Joshi

# Not a trivial task...

I eat sushi with tuna.



Aravind Joshi



# Not a trivial task...

I eat sushi with tuna.

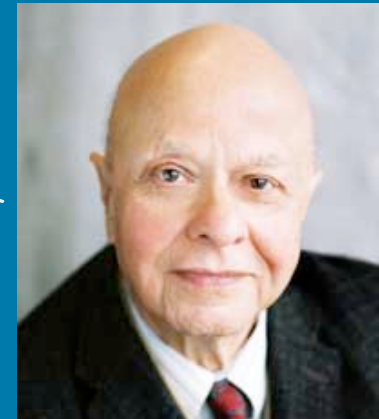


Aravind Joshi



# Not a trivial task...

I eat sushi with tuna.



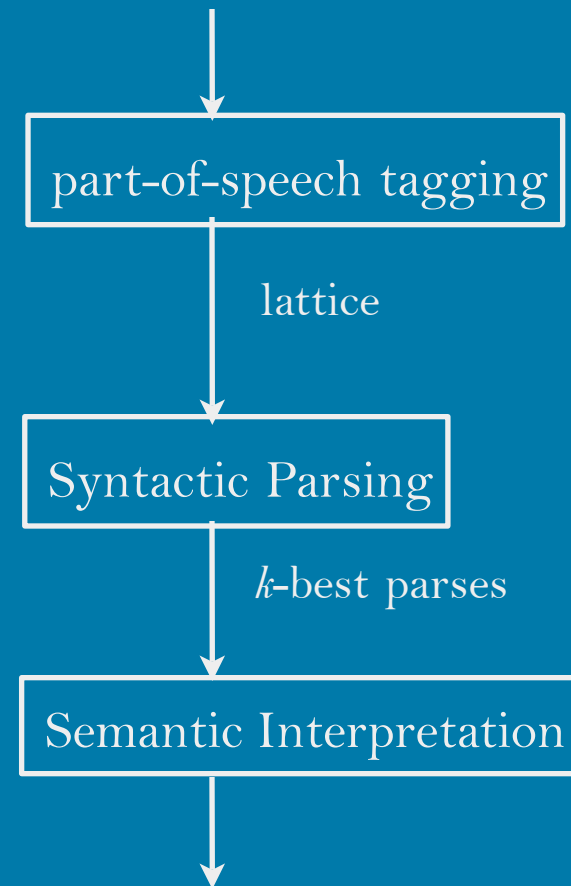
Aravind Joshi

... you'd improve  
your  $k$ -best parser



# Why $k$ -best?

- postpone disambiguation in a pipeline
  - 1-best is not always optimal in the future
  - propagate  $k$ -best lists instead of 1-best
  - e.g.: semantic role labeler uses  $k$ -best parses
- approximate the set of all possible interpretations
  - reranking (Collins, 2000)
  - minimum error training (Och, 2003)
  - online training (McDonald et al., 2005)



# In this talk...

- Formulations
  - parsing as deduction; the CKY algorithm
  - directed monotonic hypergraphs
- Algorithms
  - Algorithm 0 thru Algorithm 3
- Experiments
- Applications in Machine Translation



# Parsing as Deduction

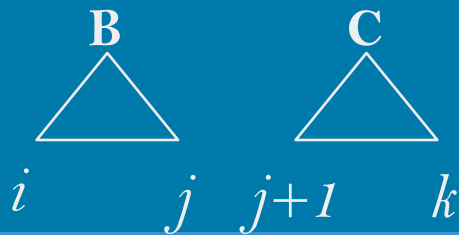
- Parsing with **context-free grammars** (CFGs)
  - Dynamic Programming (CKY algorithm)

$$\frac{(B, i, j) \quad (C, j+1, k)}{(A, i, k)} A \rightarrow B C \qquad \frac{(NP, 1, 3) \quad (VP, 4, 6)}{(S, 1, 6)} S \rightarrow NP VP$$

# Parsing as Deduction

- Parsing with **context-free grammars** (CFGs)
  - Dynamic Programming (CKY algorithm)

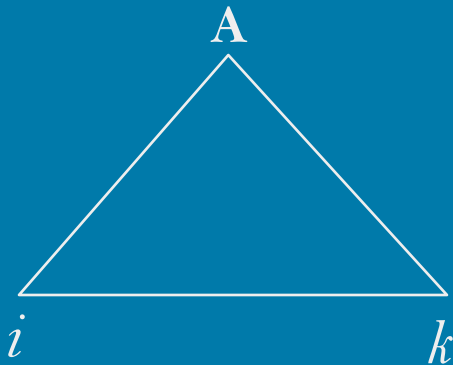
$$\frac{(B, i, j) \quad (C, j+1, k)}{(A, i, k)} A \rightarrow B C \qquad \frac{(\text{NP}, 1, 3) \quad (\text{VP}, 4, 6)}{(\text{S}, 1, 6)} \text{S} \rightarrow \text{NP VP}$$



# Parsing as Deduction

- Parsing with **context-free grammars** (CFGs)
  - Dynamic Programming (CKY algorithm)

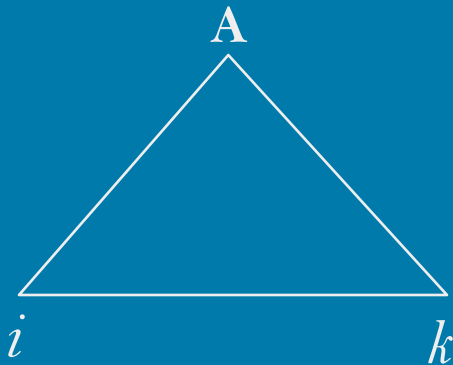
$$\frac{(B, i, j) \quad (C, j+1, k)}{(A, i, k)} A \rightarrow BC \qquad \frac{(NP, 1, 3) \quad (VP, 4, 6)}{(S, 1, 6)} S \rightarrow NP VP$$



# Parsing as Deduction

- Parsing with **context-free grammars** (CFGs)
  - Dynamic Programming (CKY algorithm)

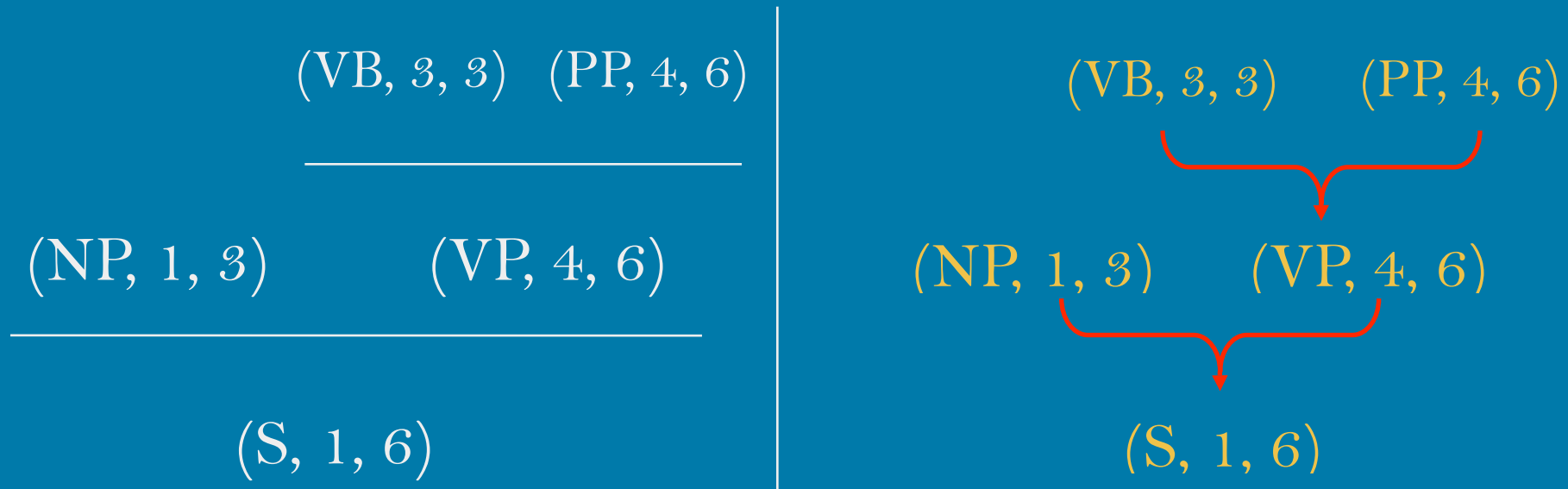
$$\frac{(B, i, j) \quad (C, j+1, k)}{(A, i, k)} A \rightarrow B C \qquad \frac{(NP, 1, 3) \quad (VP, 4, 6)}{(S, 1, 6)} S \rightarrow NP VP$$



computational complexity:  $O(n^3 |P|)$   
 $P$  is the set of productions (rules)

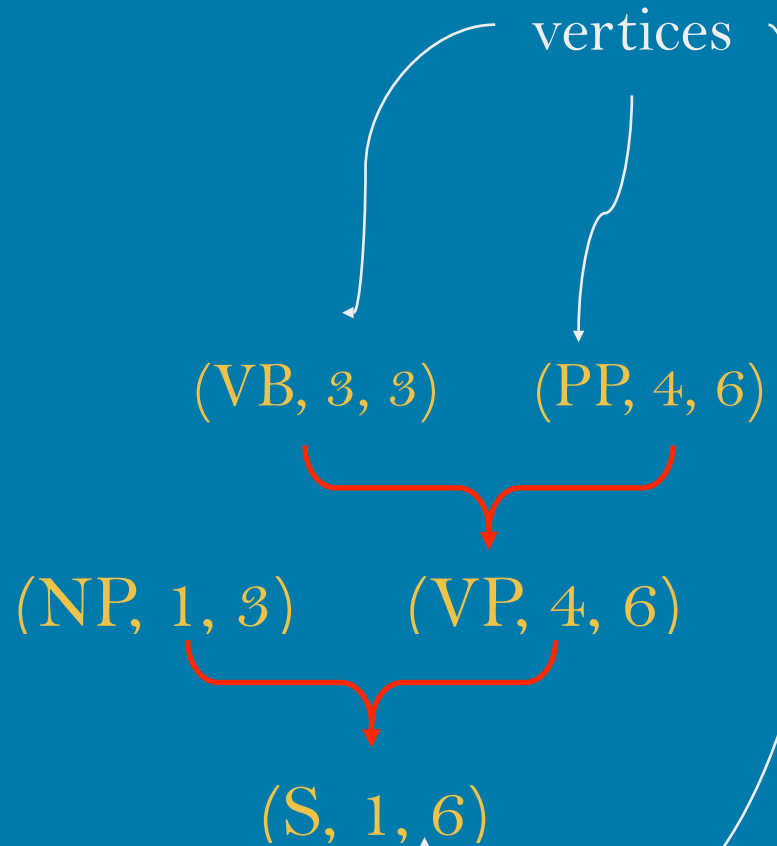
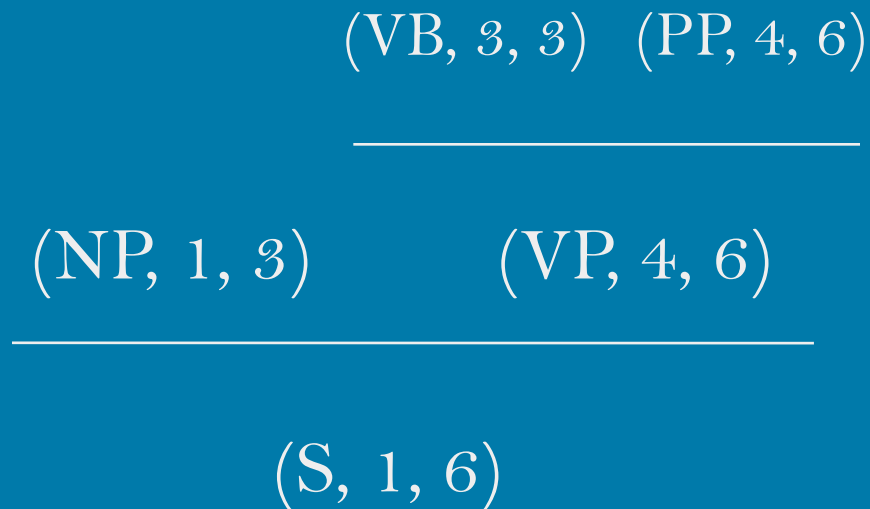
# Deduction $\Rightarrow$ Hypergraph

- hypergraph is a generalization of graph
  - each hyperedge connects several vertices to one vertex



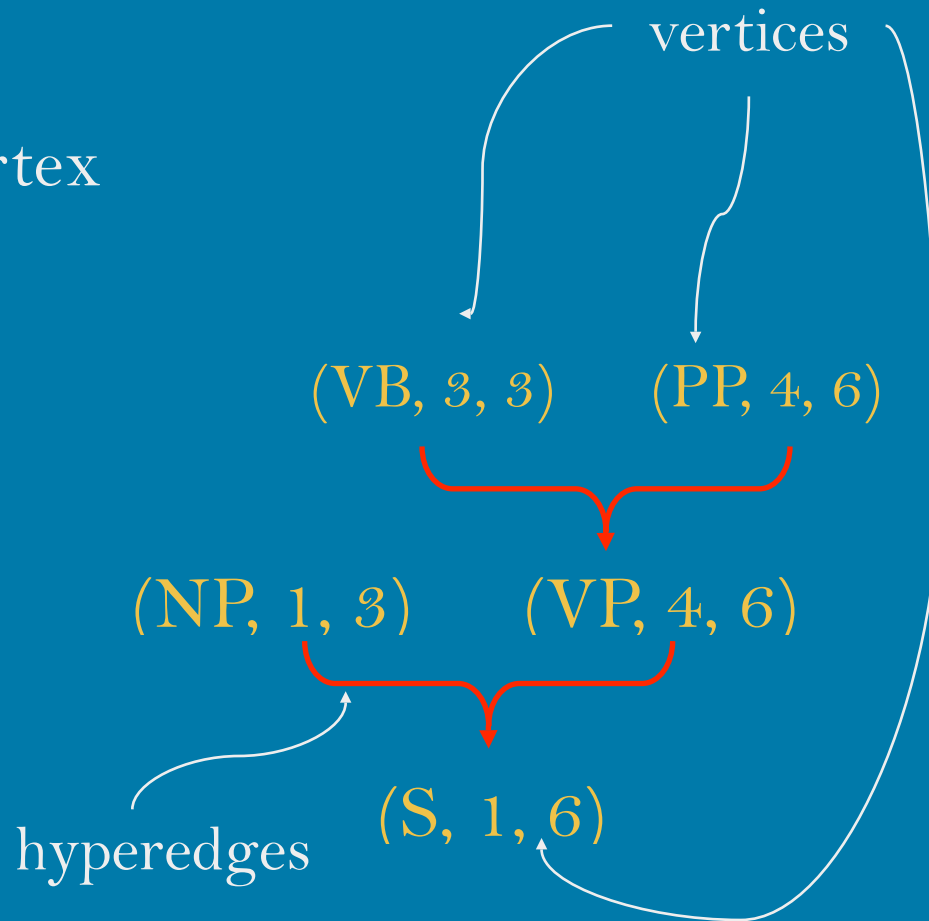
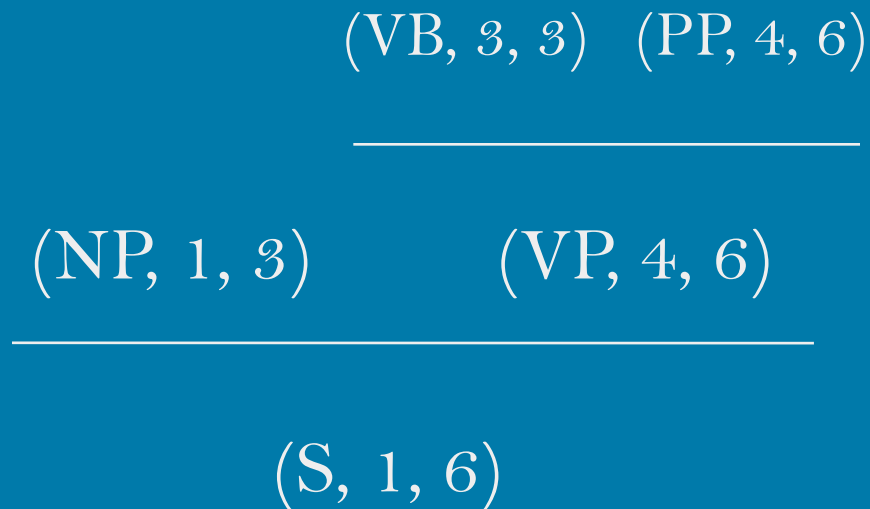
# Deduction $\Rightarrow$ Hypergraph

- hypergraph is a generalization of graph
  - each hyperedge connects several vertices to one vertex



# Deduction => Hypergraph

- hypergraph is a generalization of graph
  - each hyperedge connects several vertices to one vertex

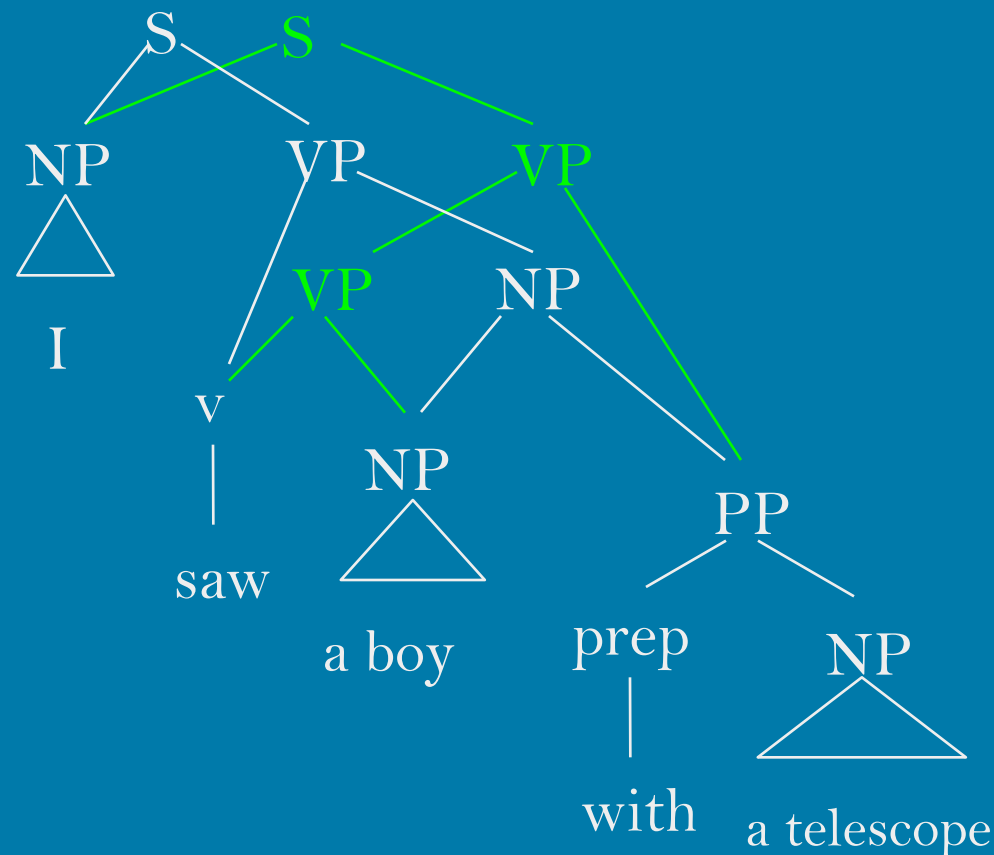


# Packed Forest as Hypergraph

packed forest

a compact representation  
of all parse trees

I saw a boy with a telescope

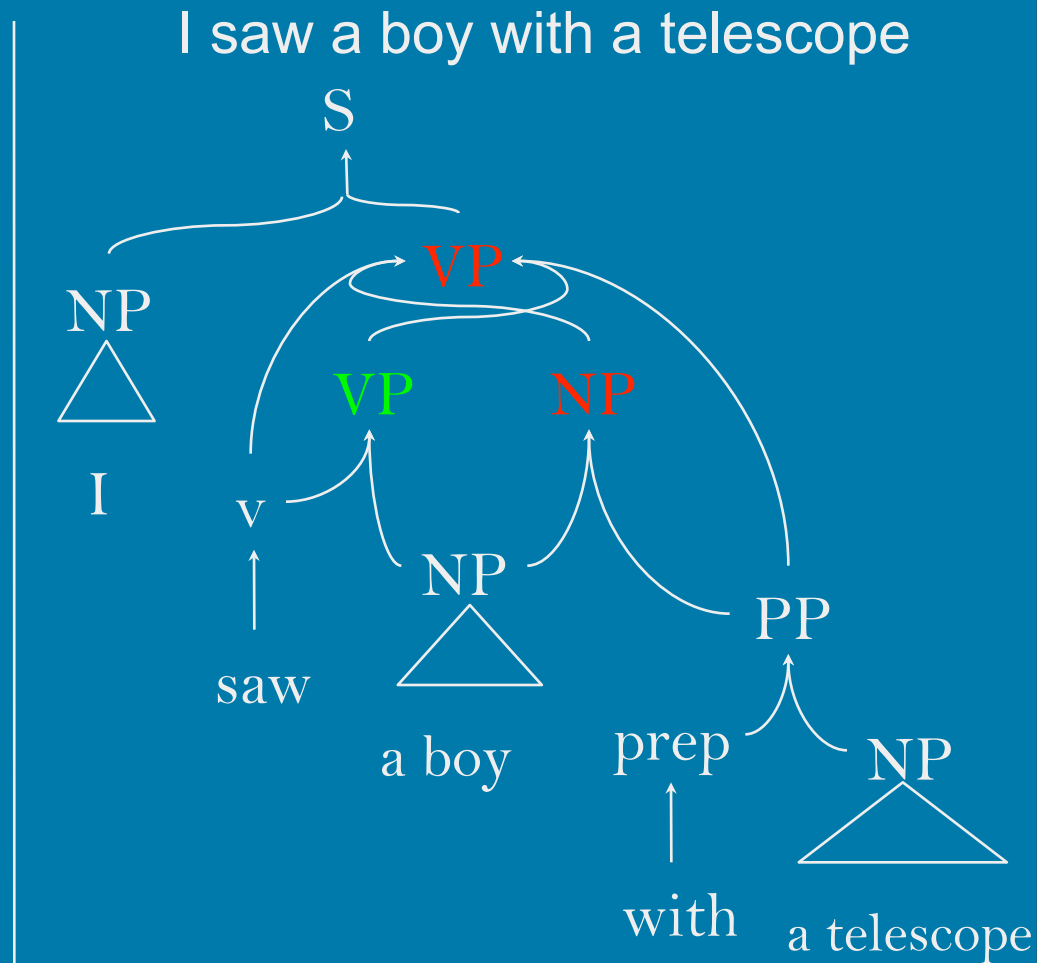




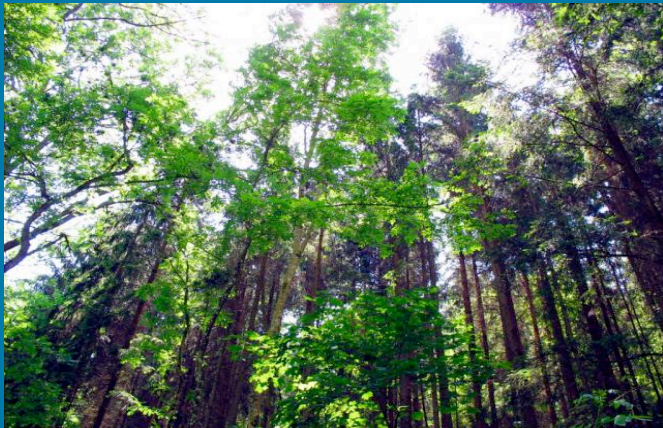
# Packed Forest as Hypergraph

packed forest

a compact representation  
of all parse trees

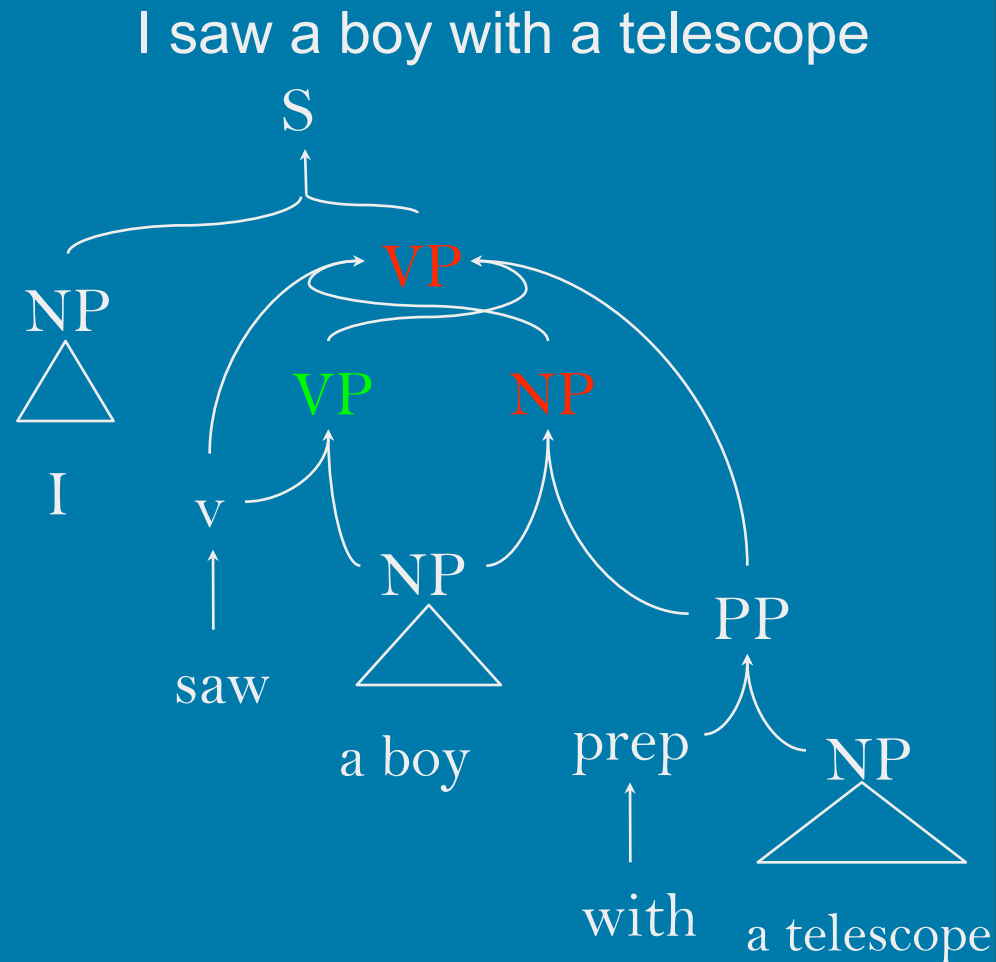


# Packed Forest as Hypergraph



packed forest

a compact representation  
of all parse trees



# Weighted Deduction/Hypergraph

$$\frac{(B, i, j): p \quad (C, j+1, k): q}{(A, i, k): f(p, q)} \quad A \longrightarrow BC$$

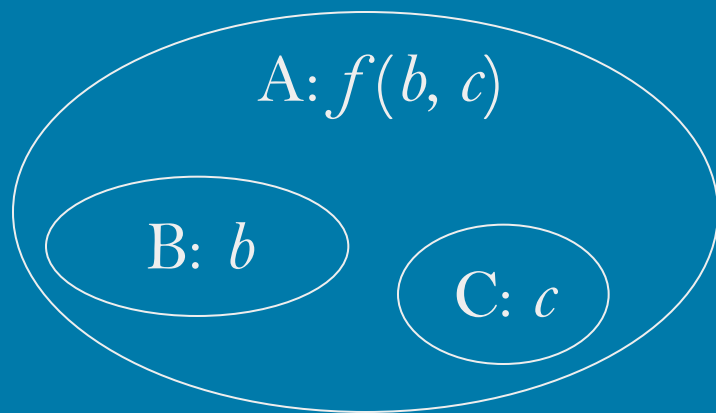
- $f$  is the weight function

e.g.: in Probabilistic Context-Free Grammars:

$$f(p, q) = p \cdot q \cdot \Pr(A \rightarrow BC)$$

# Monotonic Weight Functions

- all weight functions must be *monotonic* on each of their arguments
- optimal sub-problem property in dynamic programming



CKY example:

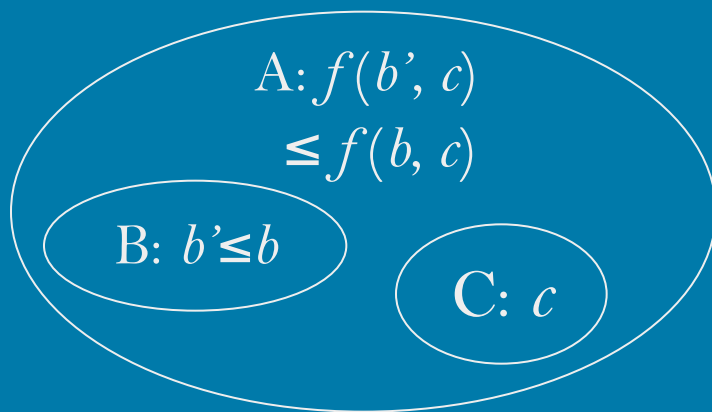
$$A = (S, 1, 5)$$

$$B = (NP, 1, 2), \quad C = (VP, 3, 5)$$

$$f(b, c) = b \cdot c \cdot \Pr(S \rightarrow NP VP)$$

# Monotonic Weight Functions

- all weight functions must be *monotonic* on each of their arguments
- optimal sub-problem property in dynamic programming



CKY example:

$$A = (S, 1, 5)$$

$$B = (NP, 1, 2), \quad C = (VP, 3, 5)$$

$$f(b, c) = b \cdot c \cdot \Pr(S \rightarrow NP VP)$$

# $k$ -best Problem in Hypergraphs

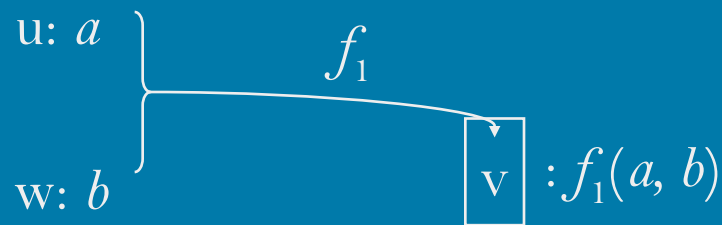
- 1-best problem
    - find the best derivation of the target vertex  $t$
  - $k$ -best problem
    - find the top  $k$  derivations of the target vertex  $t$
- in CKY,  $t = (S, 1, n)$
- assumption
    - acyclic: so that we can use topological order

# Outline

- Formulations
- **Algorithms**
  - Generic 1-best Viterbi Algorithm
  - **Algorithm 0**: naïve
  - **Algorithm 1**: hyperedge-level
  - **Algorithm 2**: vertex (item)-level
  - **Algorithm 3**: lazy algorithm
- Experiments
- Applications to Machine Translation

# Generic 1-best Viterbi Algorithm

- traverse the hypergraph in topological order (“bottom-up”)
  - for each vertex
    - for each incoming hyperedge
      - compute the result of the  $f$  function along the hyperedge
      - update the 1-best value for the current vertex if possible





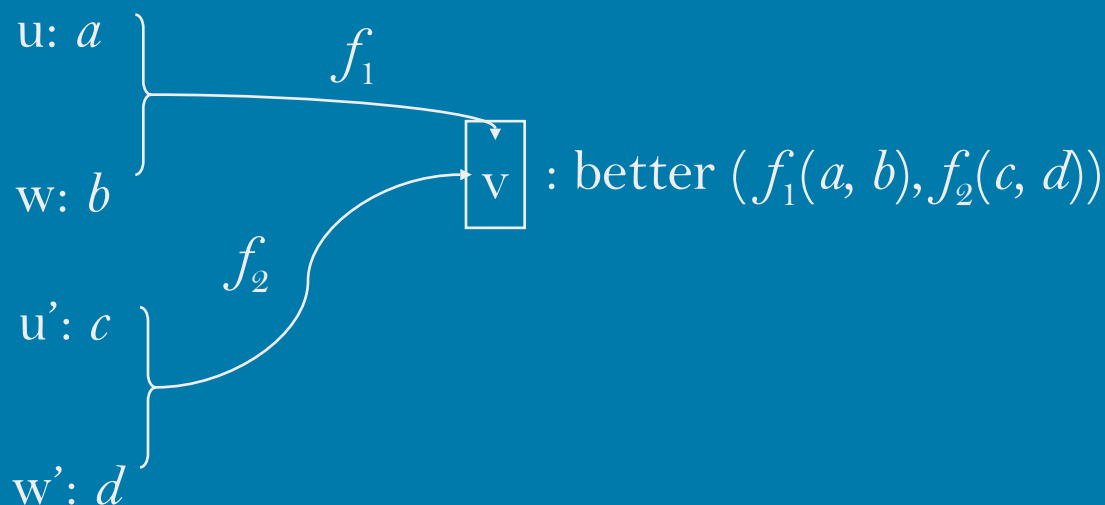
# Generic 1-best Viterbi Algorithm

- traverse the hypergraph in topological order (“bottom-up”)
  - for each vertex
    - for each incoming hyperedge
      - compute the result of the  $f$  function along the hyperedge
      - update the 1-best value for the current vertex if possible



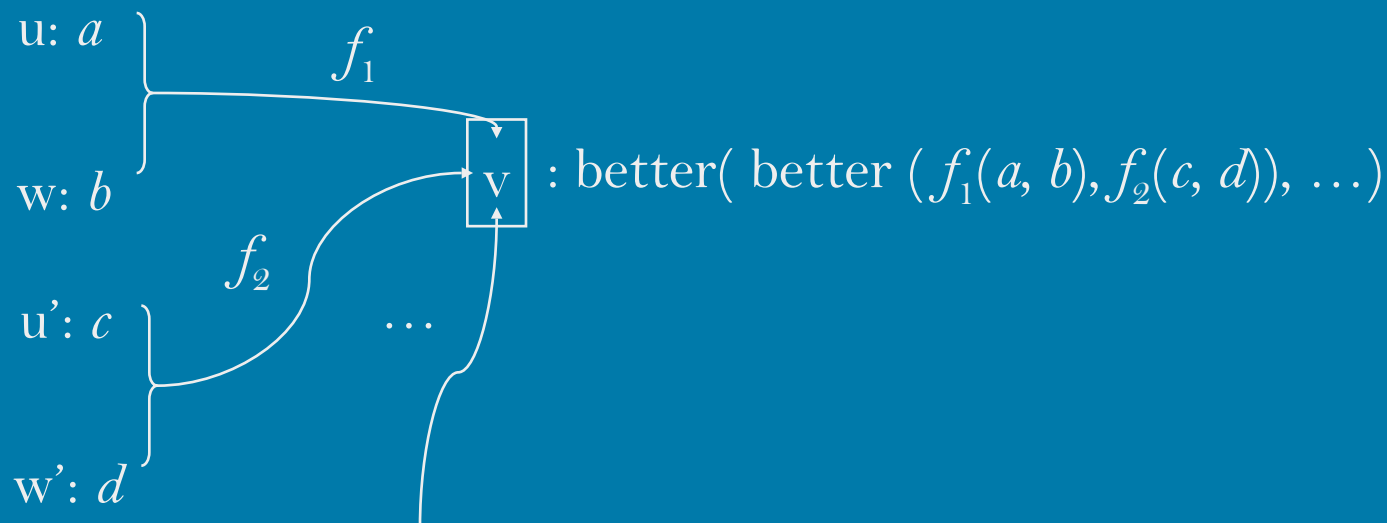
# Generic 1-best Viterbi Algorithm

- traverse the hypergraph in topological order
  - for each vertex (“bottom-up”)
    - for each incoming hyperedge
      - compute the result of the  $f$  function along the hyperedge
      - update the 1-best value for the current vertex if possible



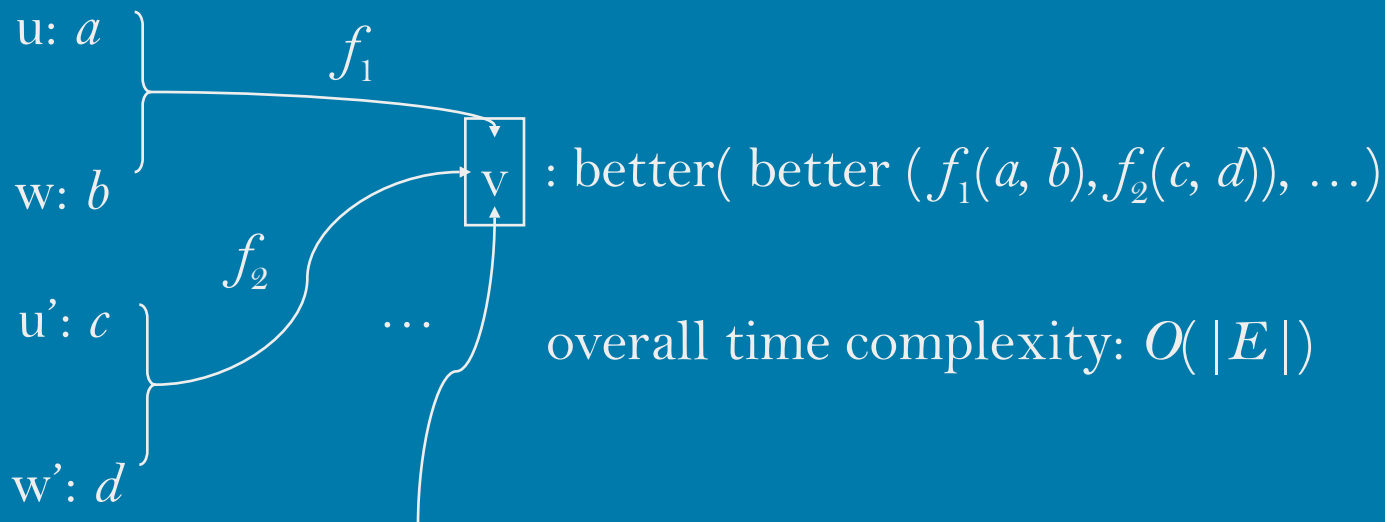
# Generic 1-best Viterbi Algorithm

- traverse the hypergraph in topological order
  - for each incoming hyperedge
    - compute the result of the  $f$  function along the hyperedge
    - update the 1-best value for the current vertex if possible



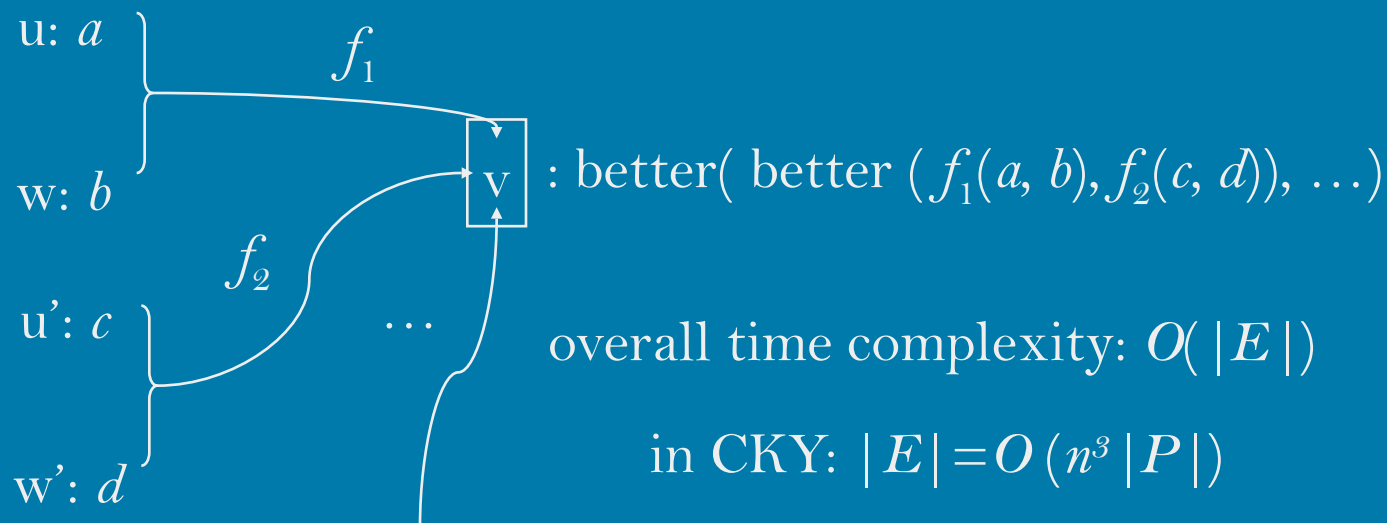
# Generic 1-best Viterbi Algorithm

- traverse the hypergraph in topological order
  - for each incoming hyperedge
    - compute the result of the  $f$  function along the hyperedge
    - update the 1-best value for the current vertex if possible

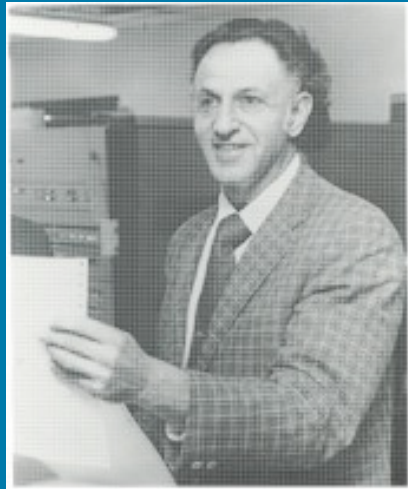


# Generic 1-best Viterbi Algorithm

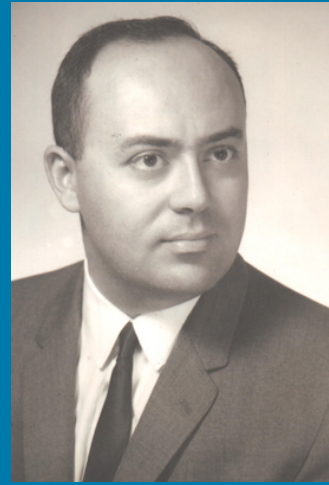
- traverse the hypergraph in topological order
  - for each incoming hyperedge
    - compute the result of the  $f$  function along the hyperedge
    - update the 1-best value for the current vertex if possible



# Dynamic Programming: 1950's

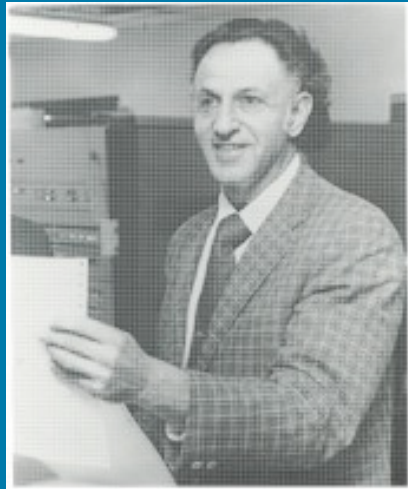


Richard Bellman

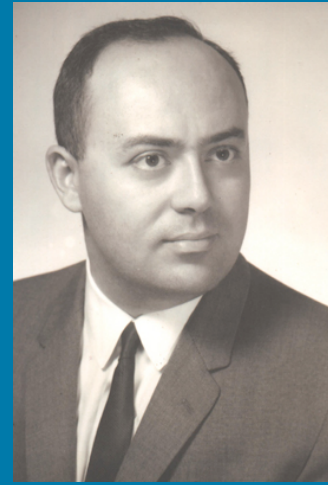


Andrew Viterbi

# Dynamic Programming: 1950's



Richard Bellman



Andrew Viterbi

We knew everything so far  
in your talk 40 years ago

# $k$ -best Viterbi algorithm 0: naïve

- straightforward  $k$ -best extension:
  - a vector of length  $k$  instead of a single value
  - vector components maintain *sorted*
  - now what's  $f(\mathbf{a}, \mathbf{b})$ ?
    - $k^2$  values -- Cartesian Product  $f(a_i, b_j)$
    - just need top  $k$  out of the  $k^2$  values

u:  $\mathbf{a}$

w:  $\mathbf{b}$

$f_1$

$\mathbf{v}$

:  $\text{mult}_k(f_1, \mathbf{a}, \mathbf{b})$

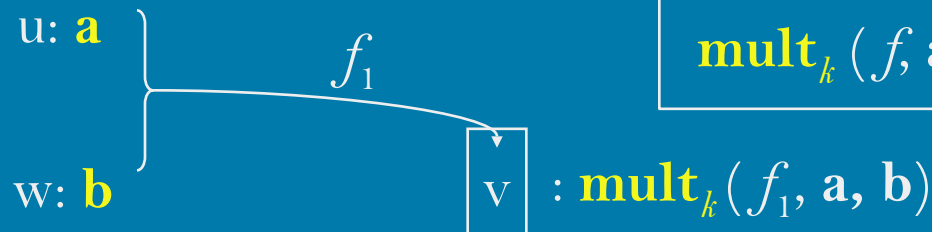
$$\text{mult}_k(f, \mathbf{a}, \mathbf{b}) = \text{top}_k \{ f(a_i, b_j) \}$$



# $k$ -best Viterbi algorithm 0: naïve

- straightforward  $k$ -best extension:
  - a vector of length  $k$  instead of a single value
  - vector components maintain *sorted*
  - now what's  $f(\mathbf{a}, \mathbf{b})$ ?
    - $k^2$  values -- Cartesian Product  $f(a_i, b_j)$
    - just need top  $k$  out of the  $k^2$  values

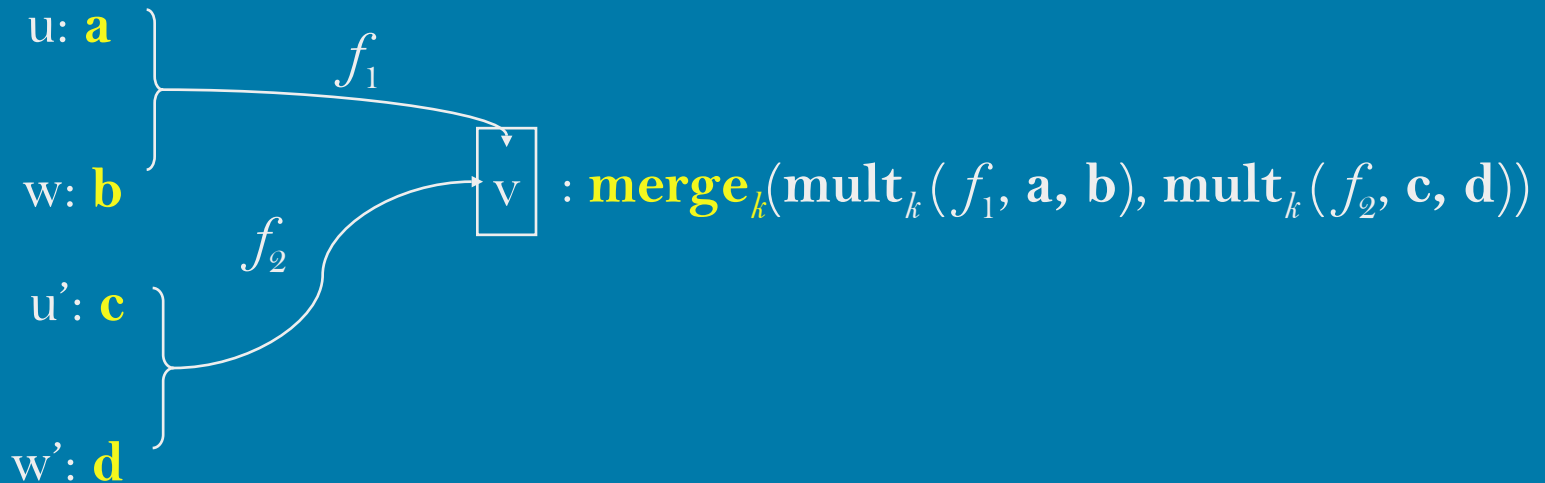
<b>b</b>	.1				
	.3				
	.4				
	.5				
		.6	.4	.3	.3
		<b>a</b>			



$$\text{mult}_k(f, \mathbf{a}, \mathbf{b}) = \text{top}_k \{ f(a_i, b_j) \}$$

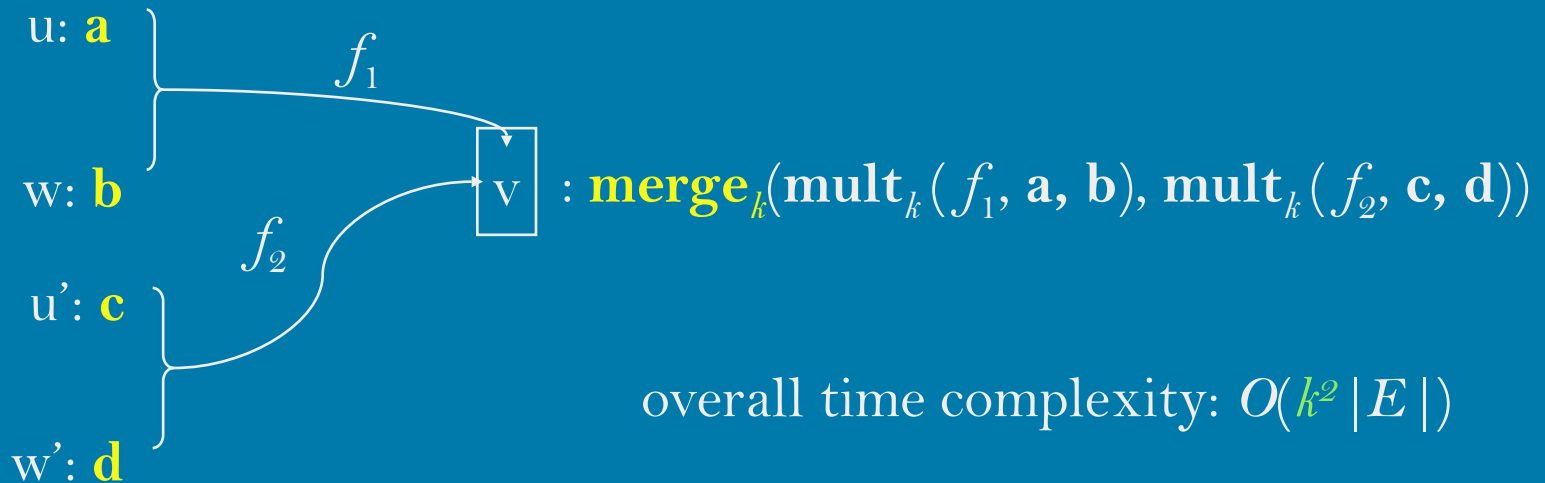
# Algorithm 0: naïve

- straightforward  $k$ -best extension:
  - a vector of length  $k$  instead of a single value
  - and how to update?
    - from two  $k$ -lengthed vectors ( $2k$  elements)
    - select the top  $k$  elements:  $O(k)$



# Algorithm 0: naïve

- straightforward  $k$ -best extension:
  - a vector of length  $k$  instead of a single value
  - and how to update?
    - from two  $k$ -lengthed vectors ( $2k$  elements)
    - select the top  $k$  elements:  $O(k)$



overall time complexity:  $O(k^2 |E|)$

# Algorithm 1: speedup $\text{mult}_k$

$$\text{mult}_k(f, \mathbf{a}, \mathbf{b}) = \text{top}_k \{ f(a_i, b_j) \}$$

	.1				
	.3				
<b>b</b>	.4				
	.5				
		.6	.4	.3	.3
					<b>a</b>

# Algorithm 1: speedup $\text{mult}_k$

$$\text{mult}_k(f, \mathbf{a}, \mathbf{b}) = \text{top}_k \{ f(a_i, b_j) \}$$

- only interested in top  $k$ , why enumerate all  $k^2$ ?

	.1				
	.3				
<b>b</b>	.4				
	.5				
		.6	.4	.3	.3
					<b>a</b>

# Algorithm 1: speedup $\text{mult}_k$

$$\text{mult}_k(f, \mathbf{a}, \mathbf{b}) = \text{top}_k \{ f(a_i, b_j) \}$$

- only interested in top  $k$ , why enumerate all  $k^2$ ?
- $\mathbf{a}$  and  $\mathbf{b}$  are sorted!

	.1				
	.3				
<b>b</b>	.4				
	.5				
		.6	.4	.3	.3
					<b>a</b>

# Algorithm 1: speedup $\text{mult}_k$

$$\text{mult}_k(f, \mathbf{a}, \mathbf{b}) = \text{top}_k \{ f(a_i, b_j) \}$$

- only interested in top  $k$ , why enumerate all  $k^2$ ?
- $\mathbf{a}$  and  $\mathbf{b}$  are sorted!
- $f$  is monotonic!

	.1				
	.3				
<b>b</b>	.4				
	.5				
		.6	.4	.3	.3
					<b>a</b>

# Algorithm 1: speedup $\text{mult}_k$

$$\text{mult}_k(f, \mathbf{a}, \mathbf{b}) = \text{top}_k \{ f(a_i, b_j) \}$$

- only interested in top  $k$ , why enumerate all  $k^2$ ?
- $\mathbf{a}$  and  $\mathbf{b}$  are sorted!
- $f$  is monotonic!
- $f(a_1, b_1)$  must be the 1-best

	.1				
	.3				
<b>b</b>	.4				
	.5				
		.6	.4	.3	.3
					<b>a</b>



# Algorithm 1: speedup $\text{mult}_k$

$$\text{mult}_k(f, \mathbf{a}, \mathbf{b}) = \text{top}_k \{ f(a_i, b_j) \}$$

- only interested in top  $k$ , why enumerate all  $k^2$ ?
- $\mathbf{a}$  and  $\mathbf{b}$  are sorted!
- $f$  is monotonic!
- $f(a_1, b_1)$  must be the 1-best
- the 2nd-best must be...
  - either  $f(a_2, b_1)$  or  $f(a_1, b_2)$

	.1				
	.3				
<b>b</b>	.4				
	.5				
		.6	.4	.3	.3
					<b>a</b>

# Algorithm 1: speedup $\text{mult}_k$

$$\text{mult}_k(f, \mathbf{a}, \mathbf{b}) = \text{top}_k \{ f(a_i, b_j) \}$$

- only interested in top  $k$ , why enumerate all  $k^2$ ?
- $\mathbf{a}$  and  $\mathbf{b}$  are sorted!
- $f$  is monotonic!
- $f(a_1, b_1)$  must be the 1-best
- the 2nd-best must be...
  - either  $f(a_2, b_1)$  or  $f(a_1, b_2)$
- what about the 3rd-best?

	.1			
	.3			
<b>b</b>	.4			
	.5			
		.6	.4	.3
				.3
				<b>a</b>

# Algorithm 1 (Demo)

$$f(a, b) = ab$$

# Algorithm 1 (Demo)

$$f(a, b) = ab$$

$b_j$	.1				
	.3				
	.4				
	.5				
		.6	.4	.3	.3
			$a_i$		

# Algorithm 1 (Demo)

$$f(a, b) = ab$$

	.1				
	.3				
$b_j$	.4				
	.5	.30			
		.6	.4	.3	.3
			$a_i$		

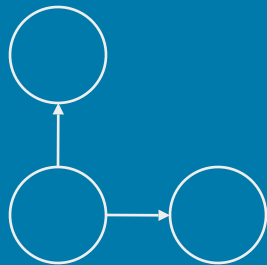
# Algorithm 1 (Demo)

$$f(a, b) = ab$$

	.1				
	.3				
$b_j$	.4	.24			
	.5	.30	.20		
		.6	.4	.3	.3
			$a_i$		

# Algorithm 1 (Demo)

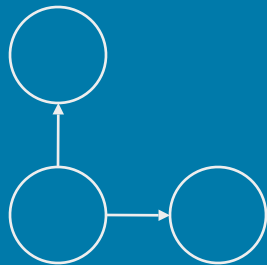
$$f(a, b) = ab$$



	.1				
	.3				
$b_j$	.4	.24			
	.5	.30	.20		
		.6	.4	.3	.3
			$a_i$		

# Algorithm 1 (Demo)

$$f(a, b) = ab$$

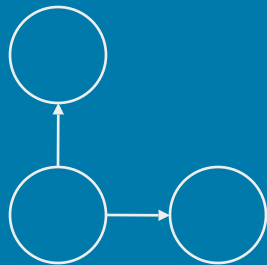


	.1				
	.3				
$b_j$	.4	.24			
	.5	.30	.20		
		.6	.4	.3	.3
			$a_i$		



# Algorithm 1 (Demo)

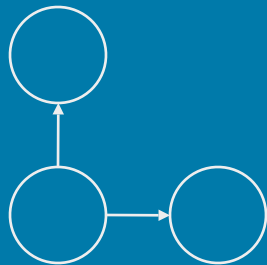
$$f(a, b) = ab$$



	.1				
	.3	.18			
$b_j$	.4	.24	.16		
	.5	.30	.20		
		.6	.4	.3	.3
			$a_i$		

# Algorithm 1 (Demo)

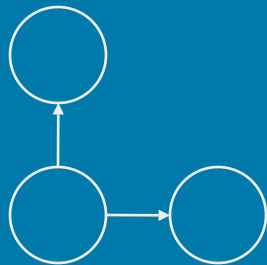
use a priority queue (heap) to store the candidates (*frontier*)



	.1				
	.3	.18			
$b_j$	.4	.24	.16		
	.5	.30	.20		
		.6	.4	.3	.3
			$a_i$		

# Algorithm 1 (Demo)

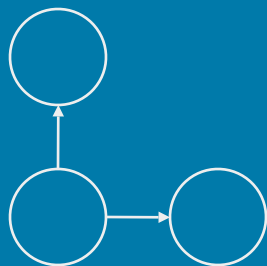
use a priority queue (heap) to store the candidates (*frontier*)



	.1				
	.3	.18			
$b_j$	.4	.24	.16		
	.5	.30	.20		
		.6	.4	.3	.3
			$a_i$		

# Algorithm 1 (Demo)

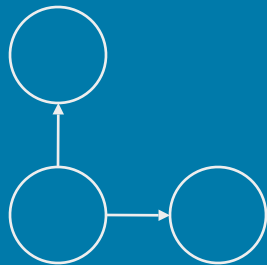
use a priority queue (heap) to store the candidates (*frontier*)



$b_j$	.1				
	.3	.18			
	.4	.24	.16		
	.5	.30	.20	.15	
		.6	.4	.3	.3
			$a_i$		

# Algorithm 1 (Demo)

use a priority queue (heap) to store the candidates (*frontier*)



	.1				
	.3	.18			
$b_j$	.4	.24	.16		
	.5	.30	.20	.15	
		.6	.4	.3	.3
			$a_i$		

in each iteration:

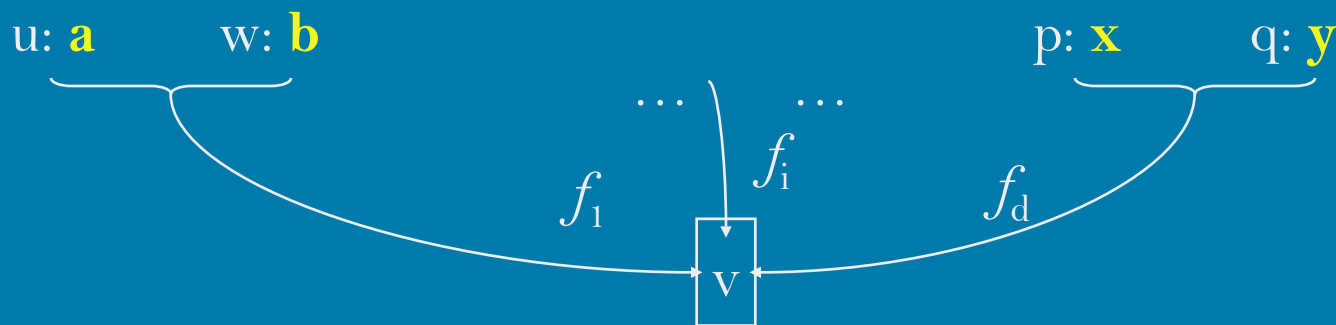
1. extract-max from the heap
2. push the two “shoulders” into the heap

$k$  iterations.

$O(k \log k |E|)$  overall time

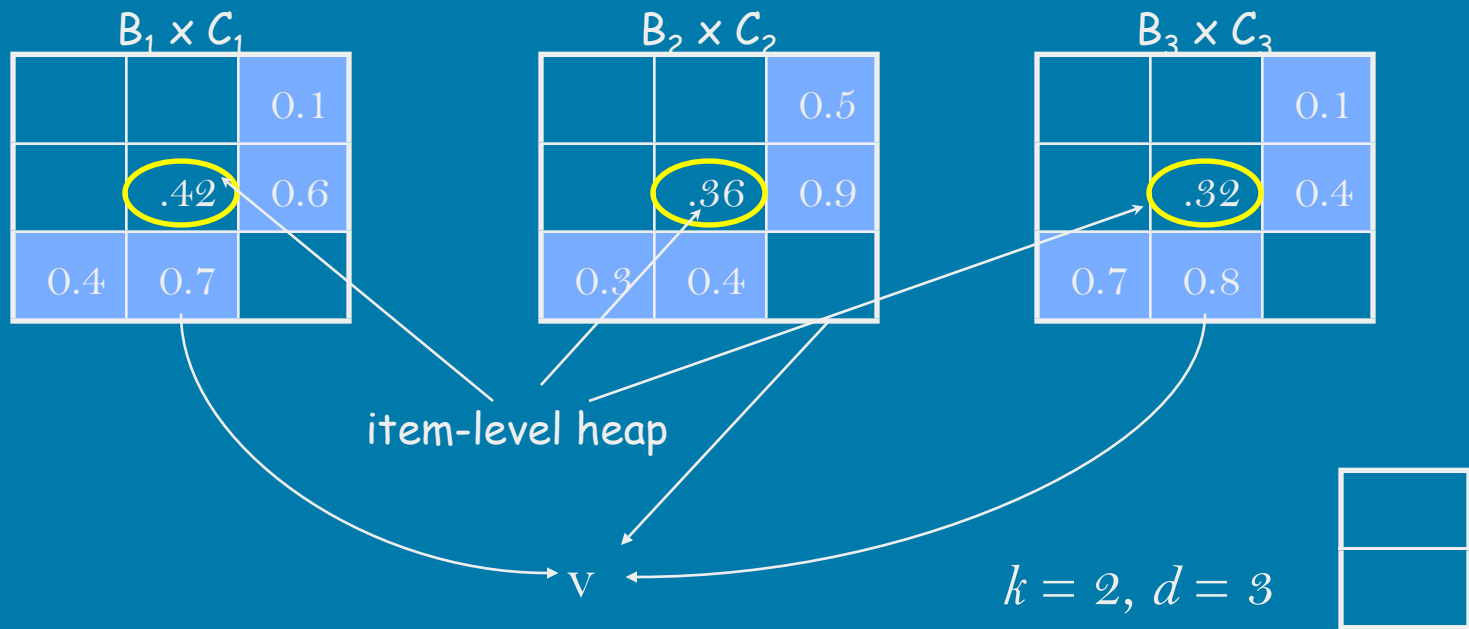
# Algorithm 2: speedup $\text{merge}_k$

- Algorithm 1 works on each hyperedge sequentially
- can we process them simultaneously?



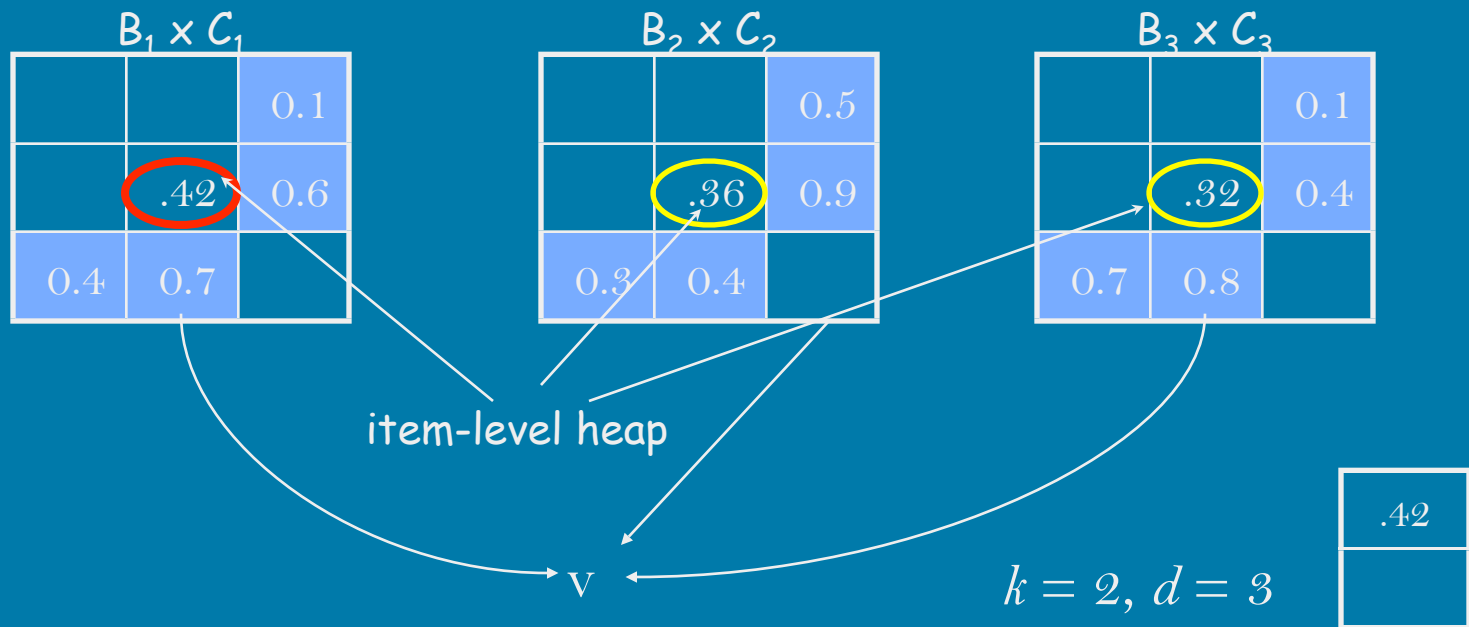
# Algorithm 2 (Demo)

starts with an initial heap of the 1-best derivations from each hyperedge



# Algorithm 2 (Demo)

pop the best (.42) and ...

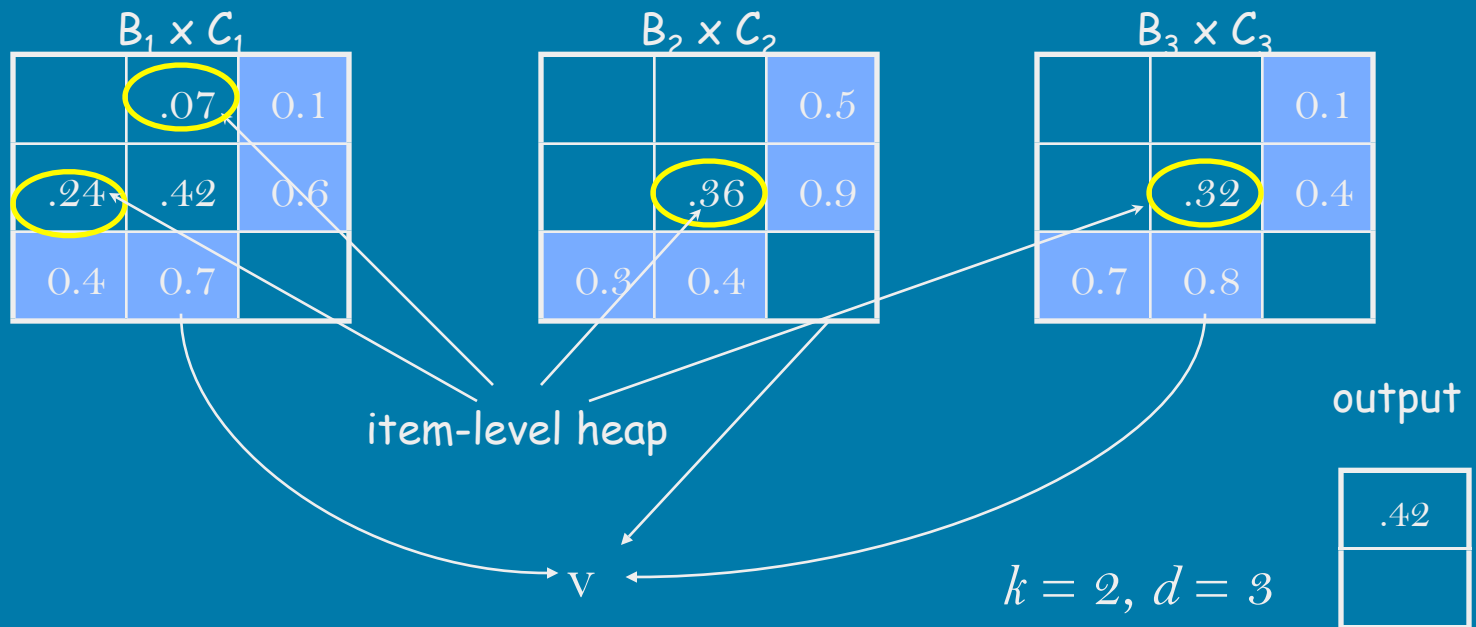




# Algorithm 2 (Demo)

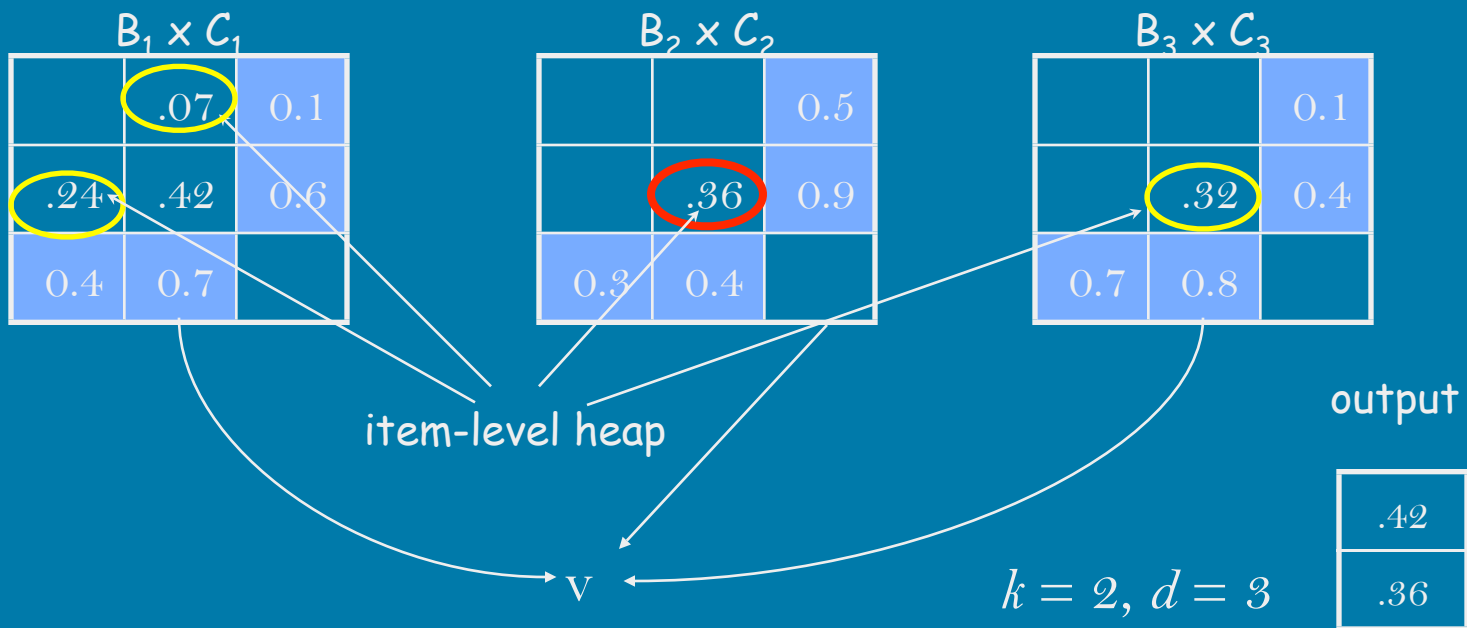
pop the best (.42) and ...

push the two successors (.07 and .24)



# Algorithm 2 (Demo)

pop the 2<sup>nd</sup>-best (.36)



# Algorithm 3: Offline (lazy)

- from Algorithm 0 to Algorithm 2:
  - delaying the calculations until needed -- lazier
  - larger locality
- even lazier... (one step further)
  - we are interested in the  $k$ -best derivations of the **final item only!**

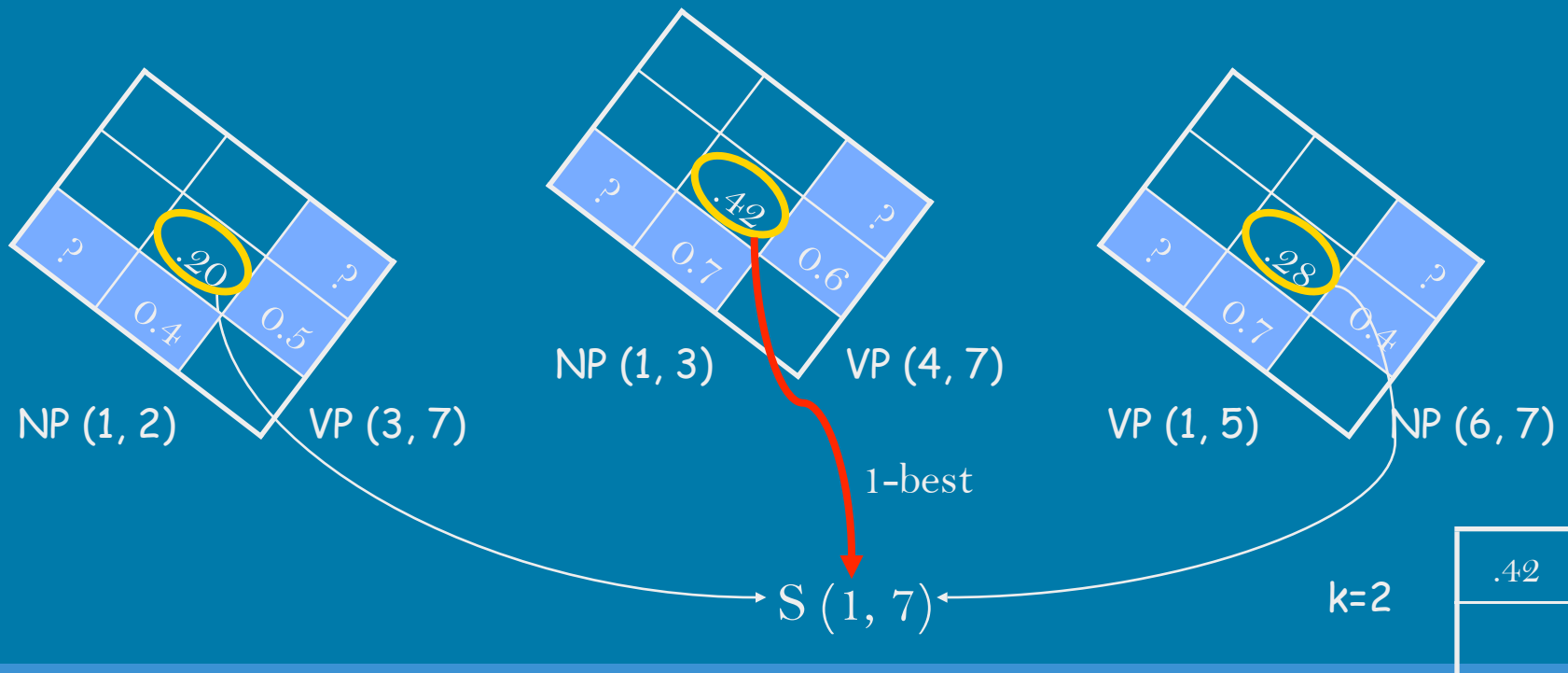
# Algorithm 3: Offline (lazy)

- forward phase
  - do a normal 1-best search till the final item
  - *but* construct the hypergraph (*forest*) along the way
- recursive backward phase
  - ask the final item: what's your 2nd-best?
  - final item will propagate this question till the leaves
  - then ask the final item: what's your 3rd-best?

# Algorithm 3 demo

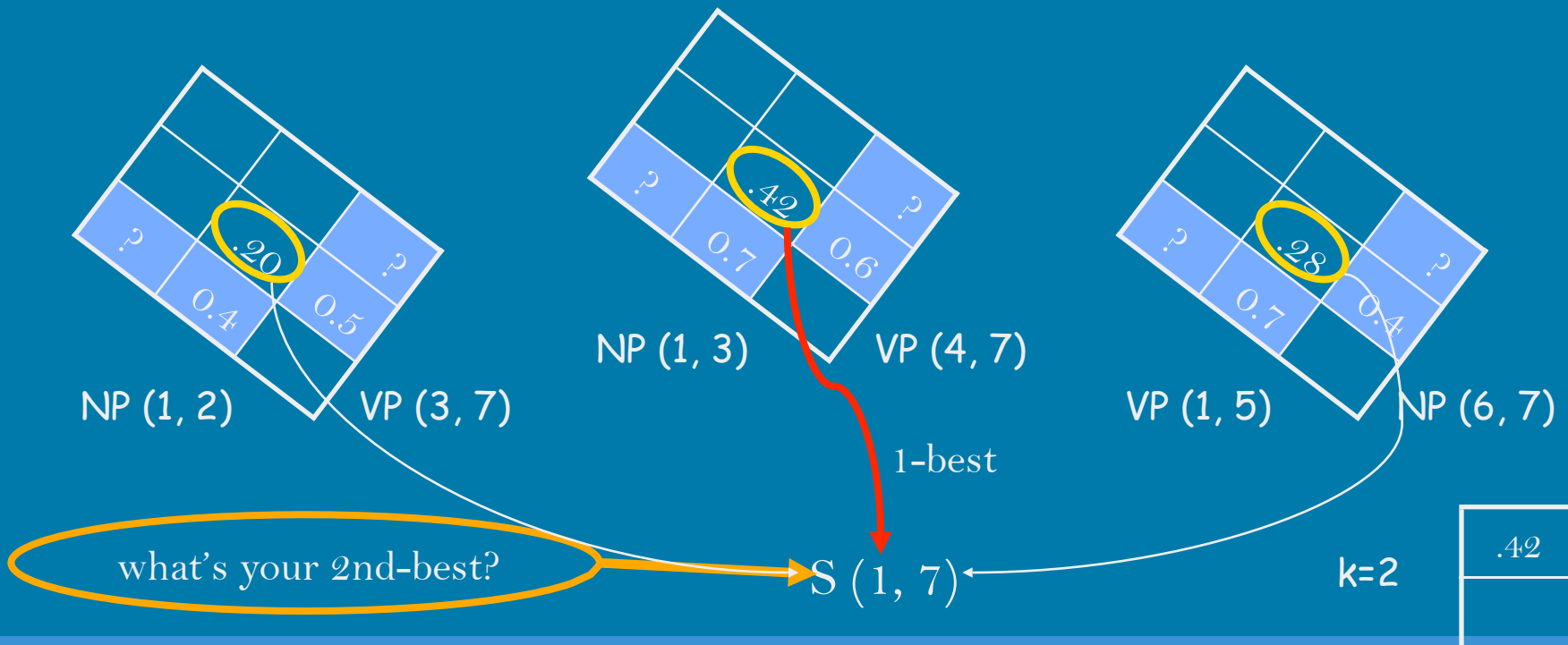
after the “forward” step (1-best parsing):

forest = 1-best derivations from each hyperedge

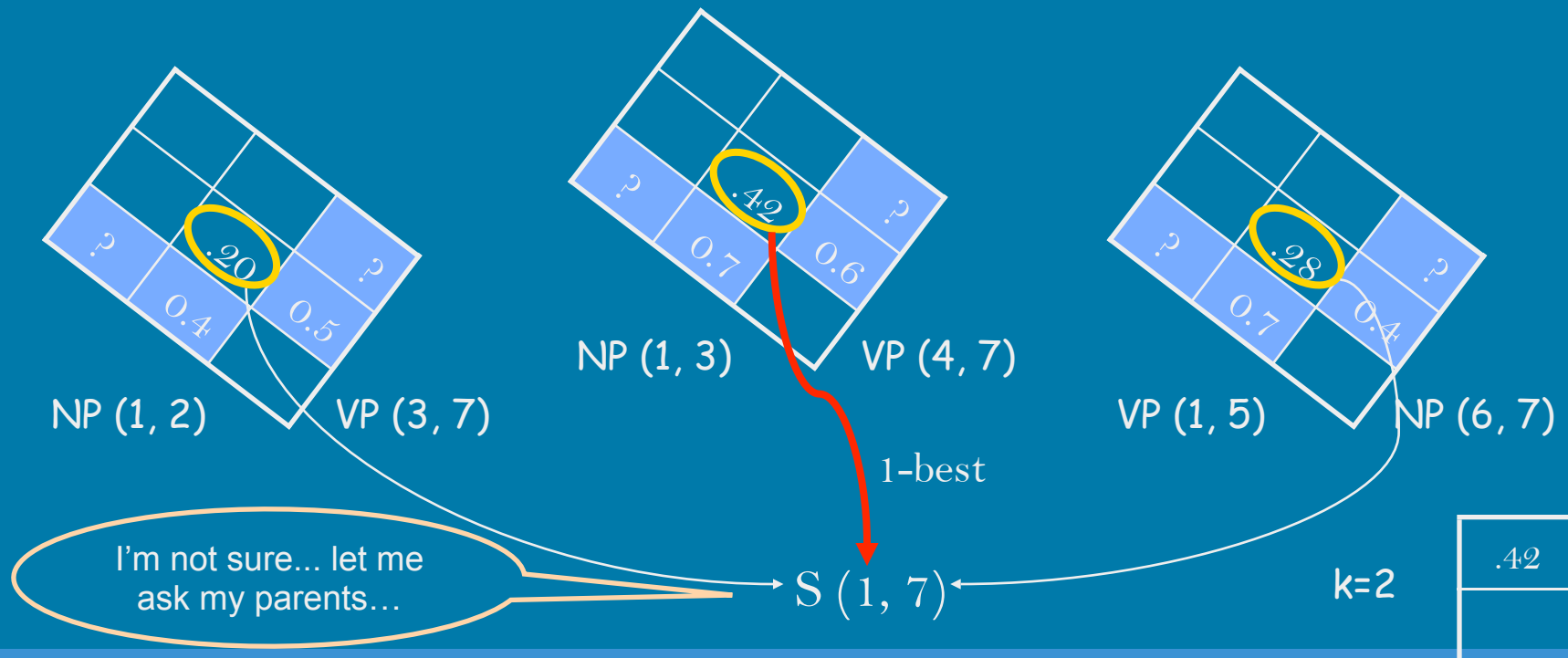


# Algorithm 3 demo

now the backward step

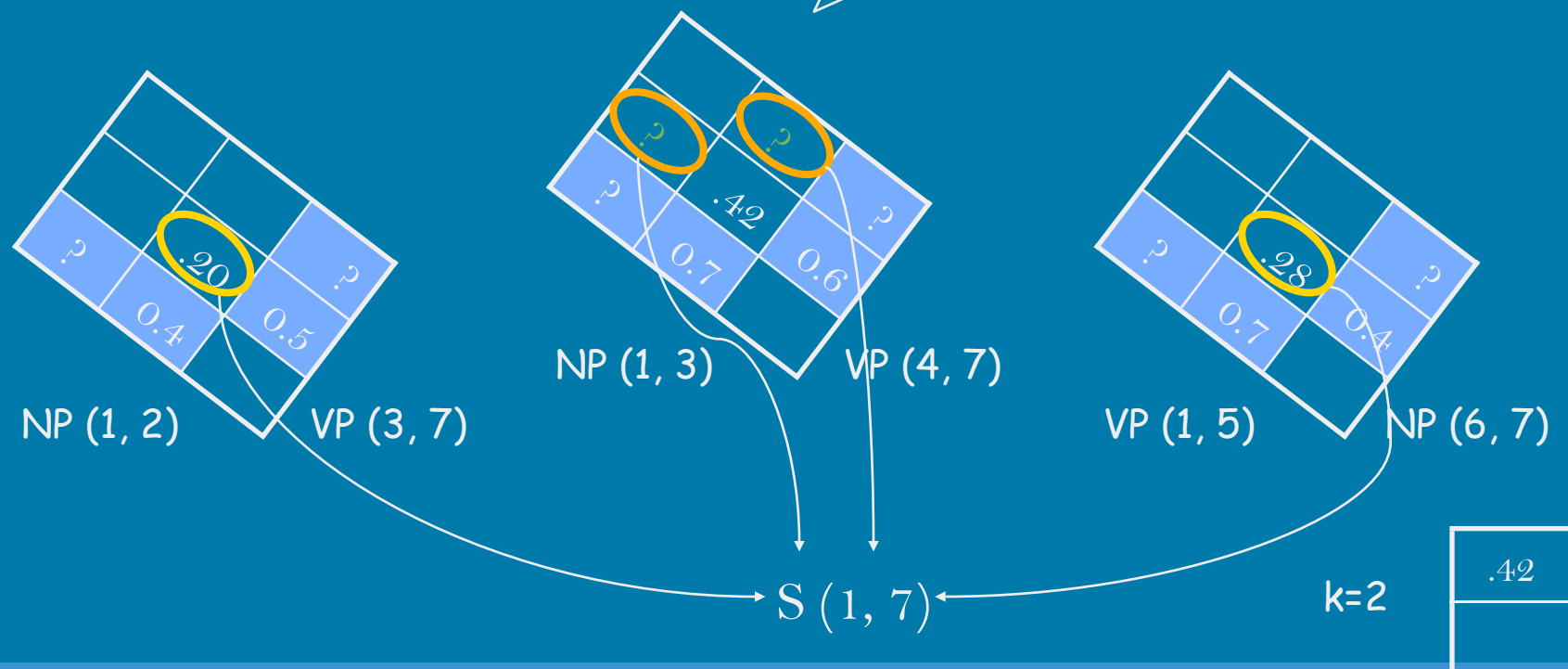


# Algorithm 3 demo



# Algorithm 3 demo

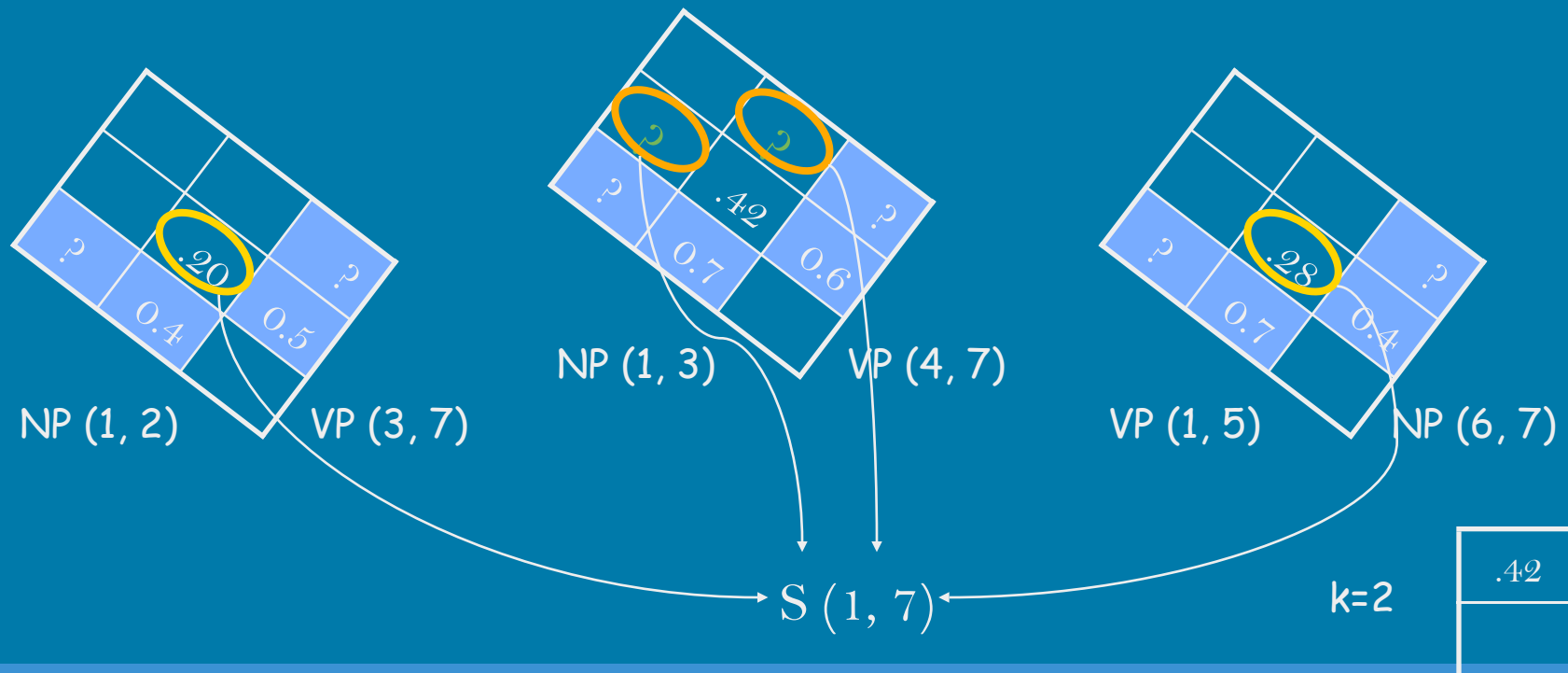
well, it must be either ... or ...





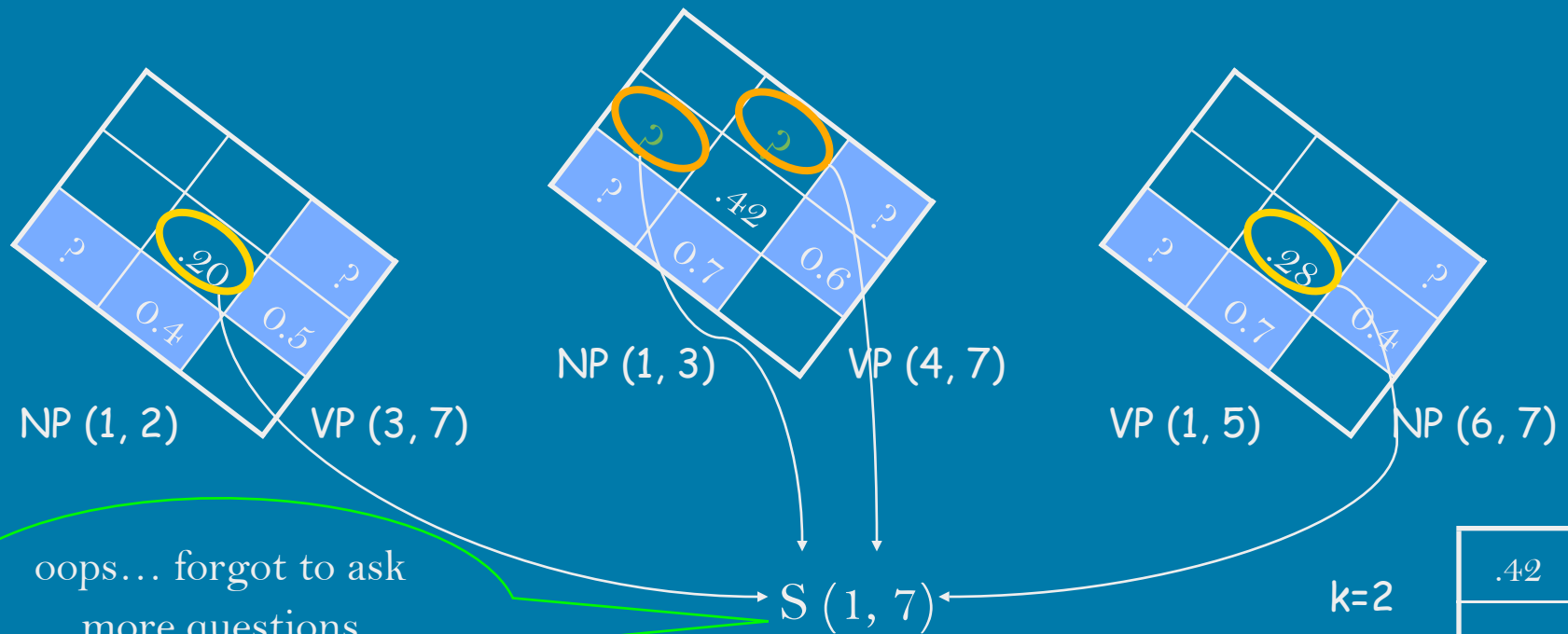
# Algorithm 3 demo

but wait a minute... did you already know the ?'s ?

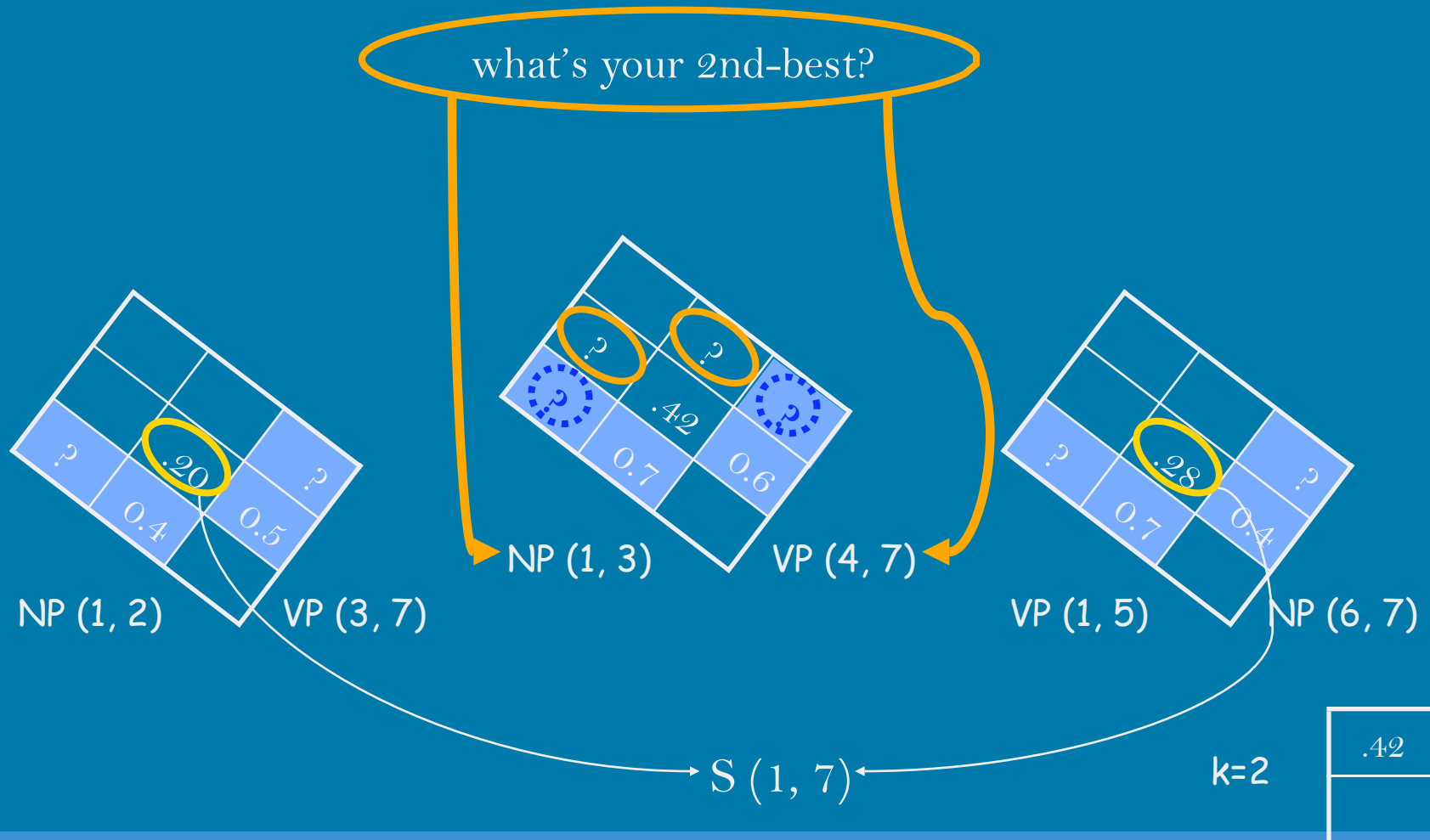


# Algorithm 3 demo

but wait a minute... did you already know the ?'s ?

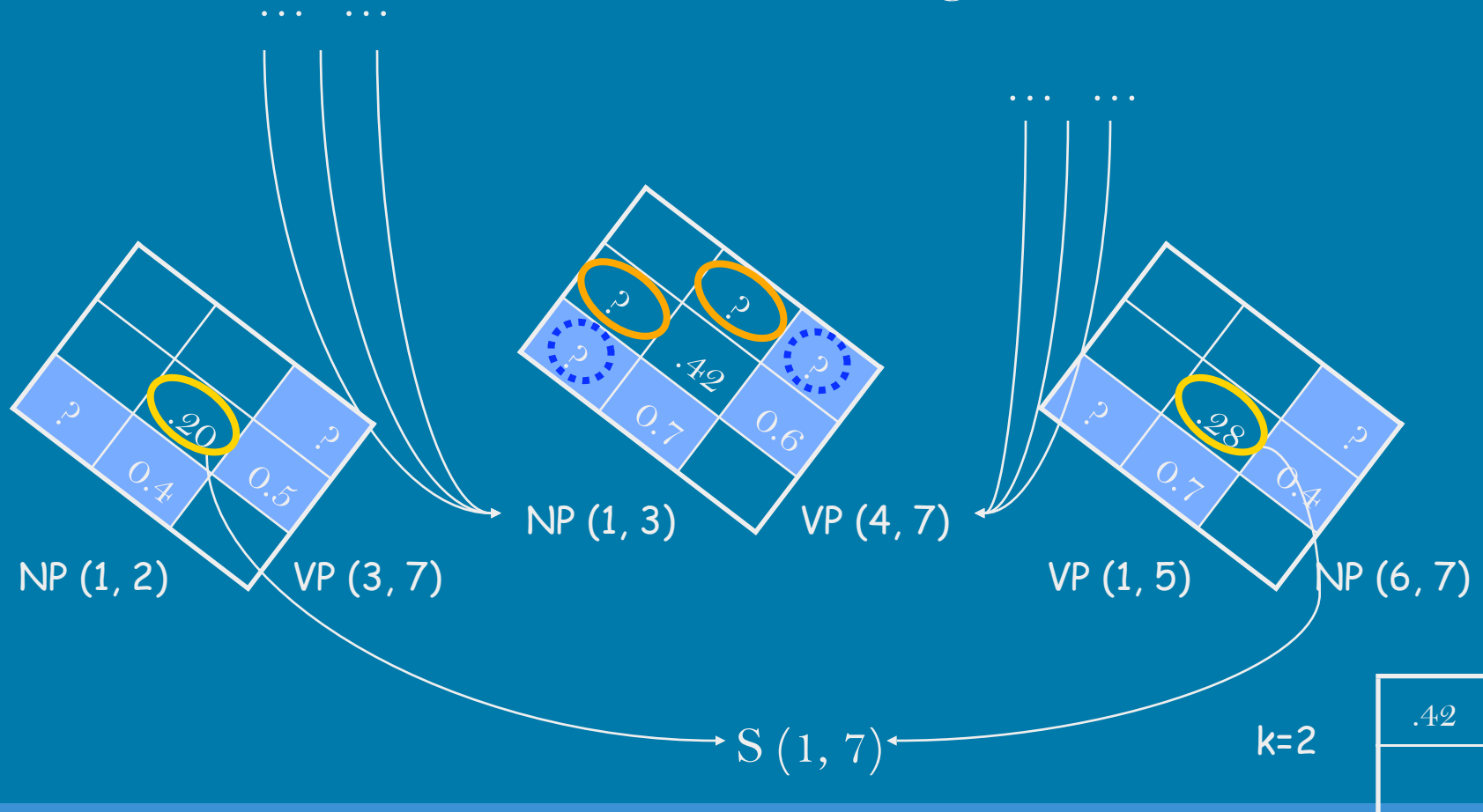


# Algorithm 3 demo



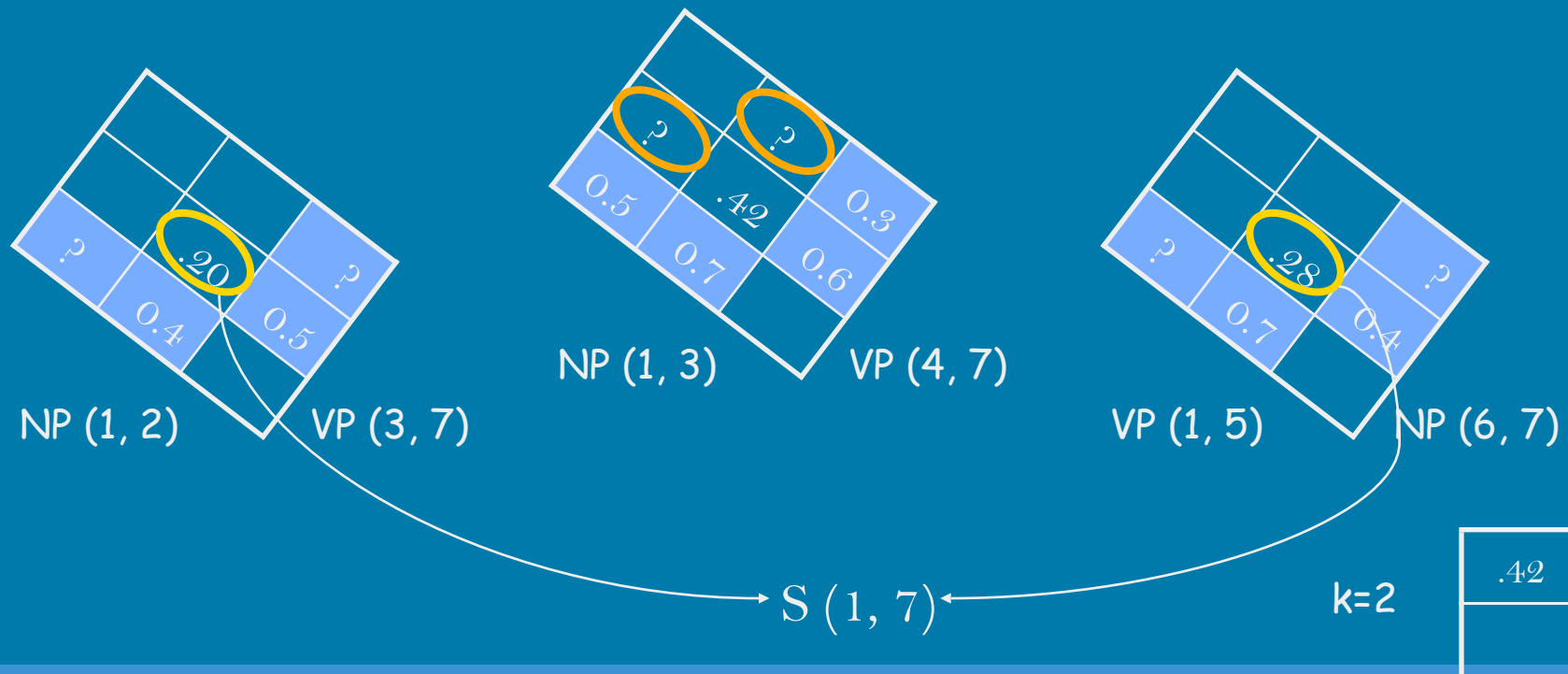
# Algorithm 3 demo

recursion goes on to the leaf nodes



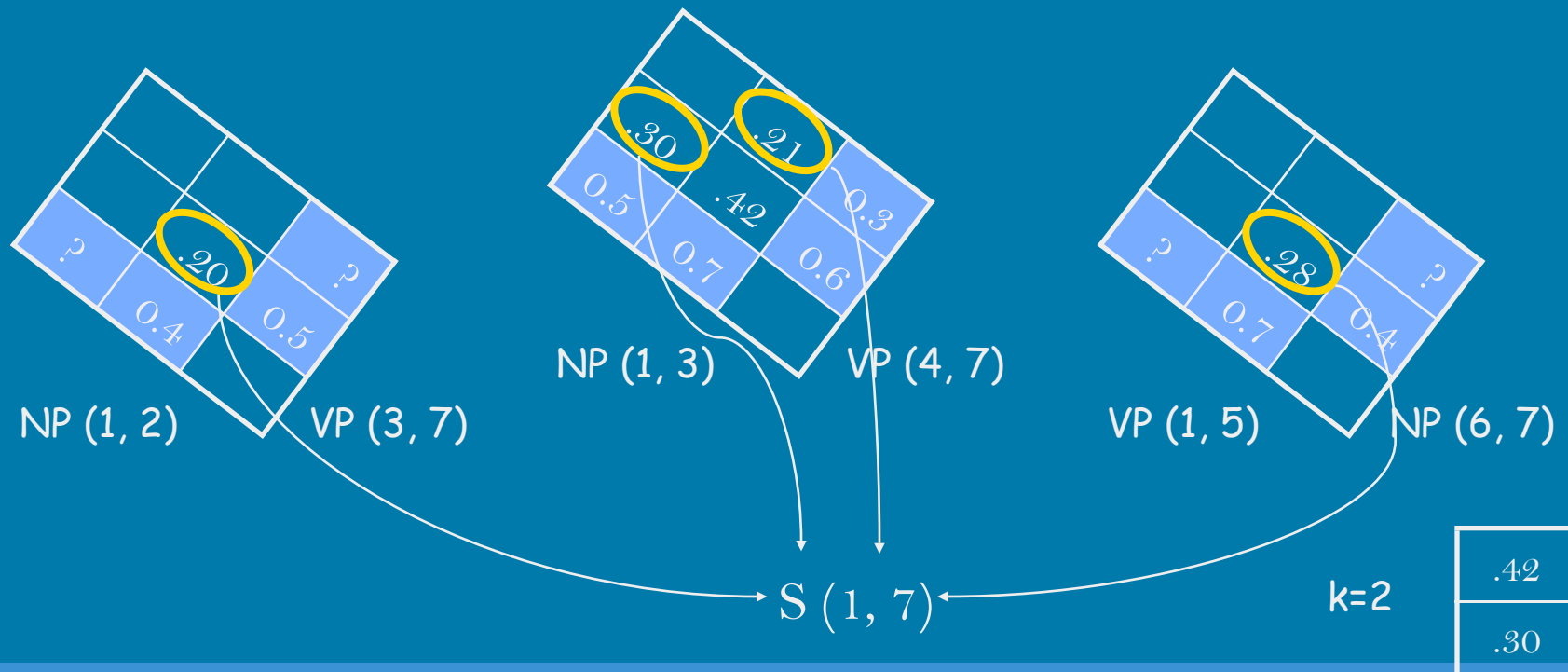
# Algorithm 3 demo

and reports back the numbers...



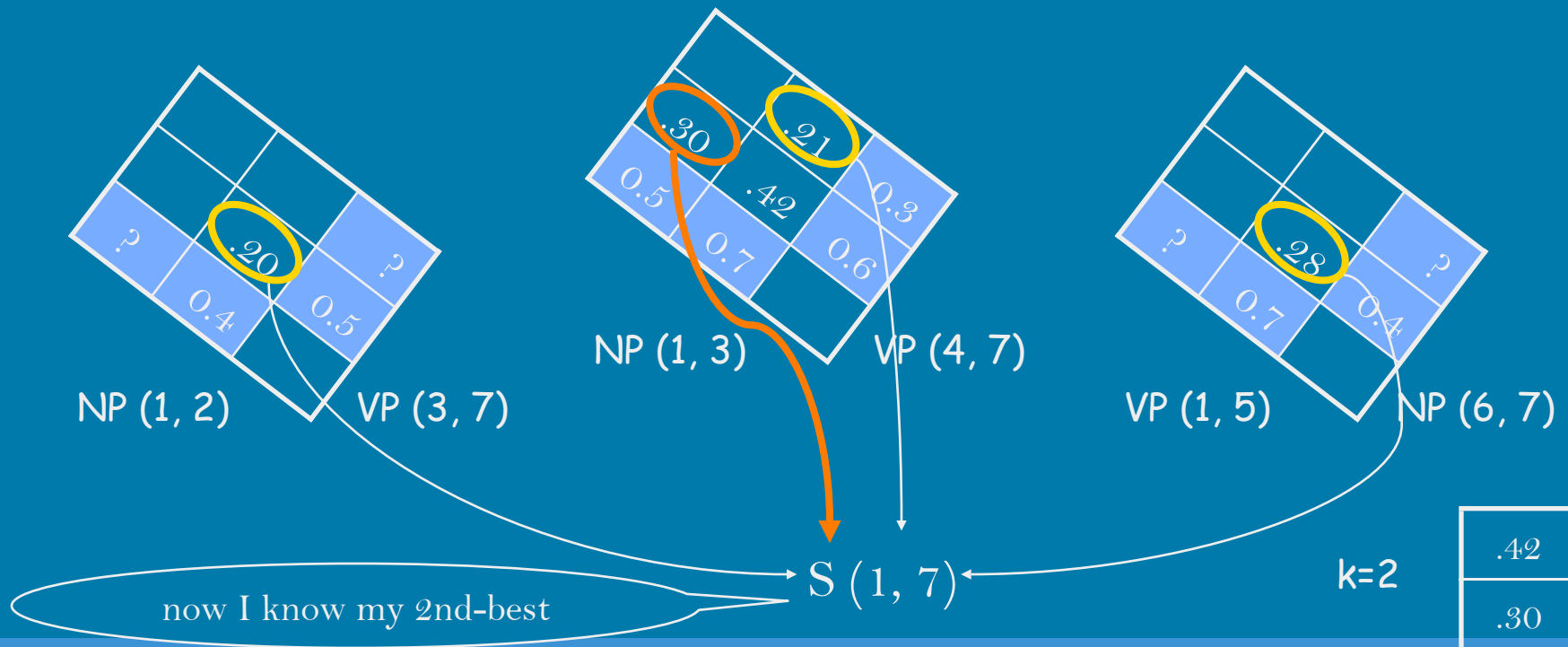
# Algorithm 3 demo

push .30 and .21 to the candidate heap (priority queue)



# Algorithm 3 demo

pop the root of the heap (.30)

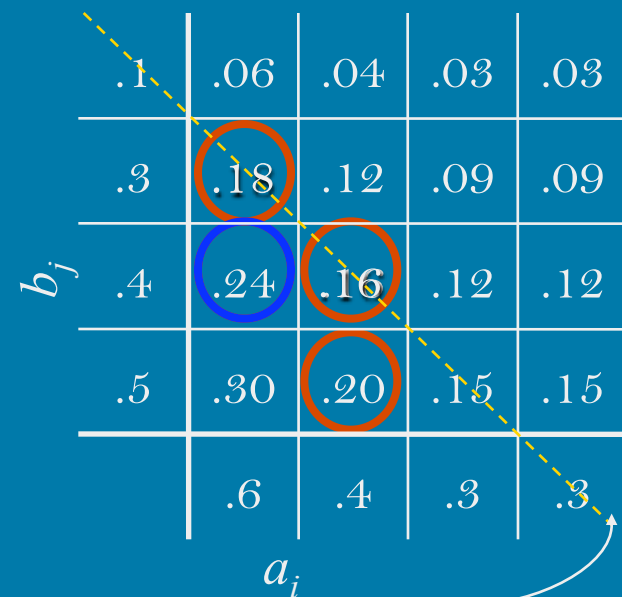


# Interesting Properties

- 1-best is best everywhere (all decisions optimal)
- 2nd-best is optimal everywhere except one decision
  - and that decision must be 2nd-best
  - and it's the best of all 2nd-best decisions
- so what about the 3rd-best?
- $k$ th-best is...

$$\sum_{\delta \in \Delta} (\text{rank}(\delta) - 1) \leq k - 1$$

local picture:



$$(i - 1) + (j - 1) = 3 - 1$$



# Summary of Algorithms

Algorithms	Time Complexity	Locality
1-best	$O( E )$	hyperedge
alg. 0: naïve	$O(k^a  E )$	hyperedge ( $\text{mult}_k$ )
alg. 1	$O(k \log k  E )$	hyperedge ( $\text{mult}_k$ )
alg. 2	$O( E  +  \mathcal{V}  k \log k)$	item ( $\text{merge}_k$ )
alg. 3: lazy	$O( E  +  \mathcal{D}  k \log k)$	global

for CKY:  $a=2$ ,  $|E| = O(n^3 |P|)$ ,  $|\mathcal{V}| = O(n^2 |N|)$ ,  $|\mathcal{D}| = O(n)$   
 $a$  is the arity of the grammar

# Outline

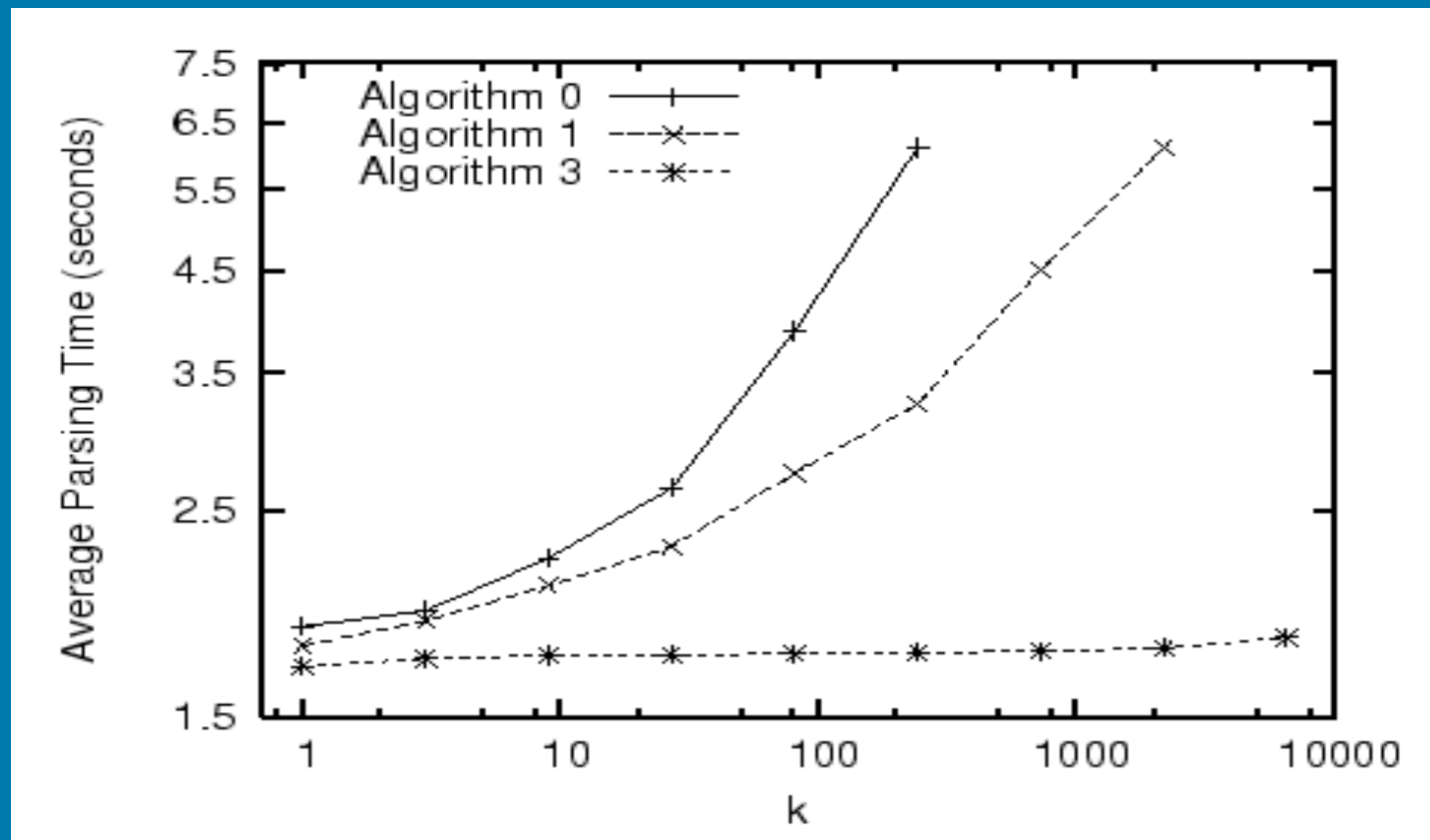
- Formulations
- Algorithms: Alg.0 thru Alg. 3
- **Experiments**
  - Collins/Bikel Parser
  - both efficiency and accuracy
- Applications in Machine Translation

# Background: Statistical Parsing

- Probabilistic Grammar
  - induced from a treebank (Penn Treebank)
- State-of-the-art Parsers
  - Collins (1999), Bikel (2004), Charniak (2000), etc.
- Evaluation of Accuracy
  - PARSEVAL: tree-similarity (English treebank: ~90%)
- Previous work on  $k$ -best Parsing:
  - Collins (2000): turn off dynamic programming
  - Charniak/Johnson (2005): coarse-to-fine, still too slow

# Efficiency

Implemented Algorithms 0, 1, 3 on top of Collins/Bikel Parser  
Average (wall-clock) time on Penn Treebank (per sentence):

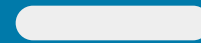


$$O(|E| + |D| k \log k)$$

# Oracle Reranking

given  $k$  parses of a sentence

gold standard  
“correct” parse



accuracy  
100%

- **oracle reranking**: pick the best parse according to the gold-standard

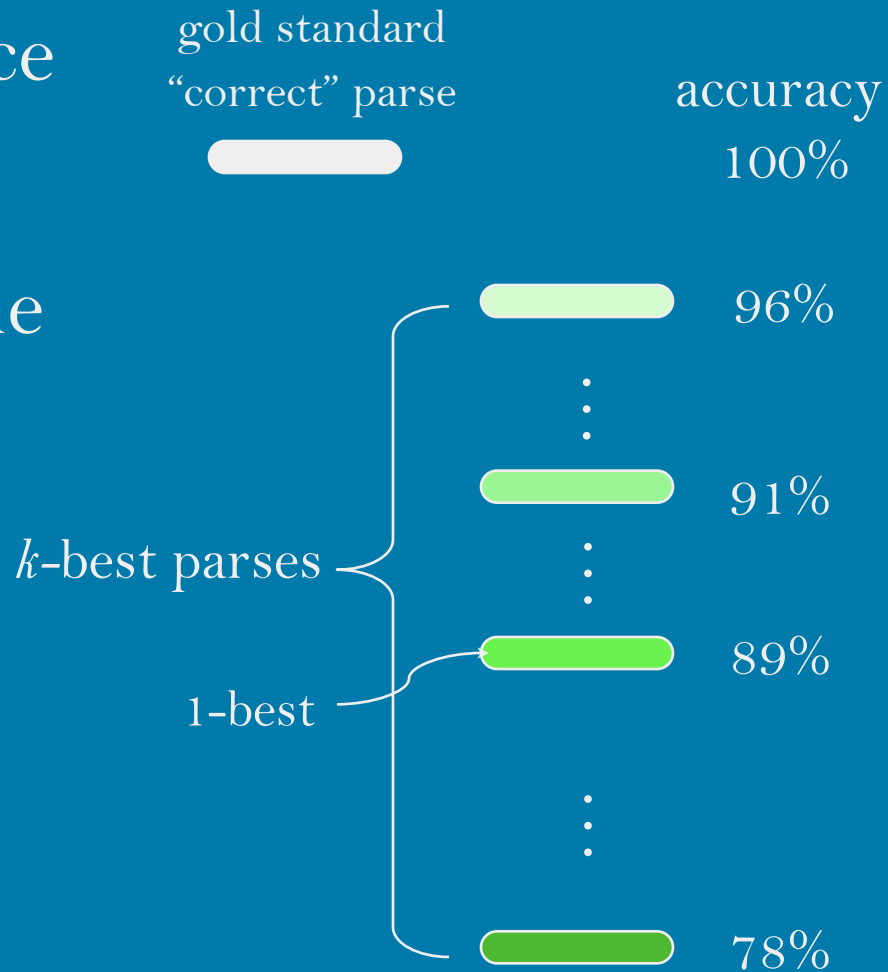
- **real reranking**: pick the best parse according to the score function

1-best  89%

# Oracle Reranking

given  $k$  parses of a sentence

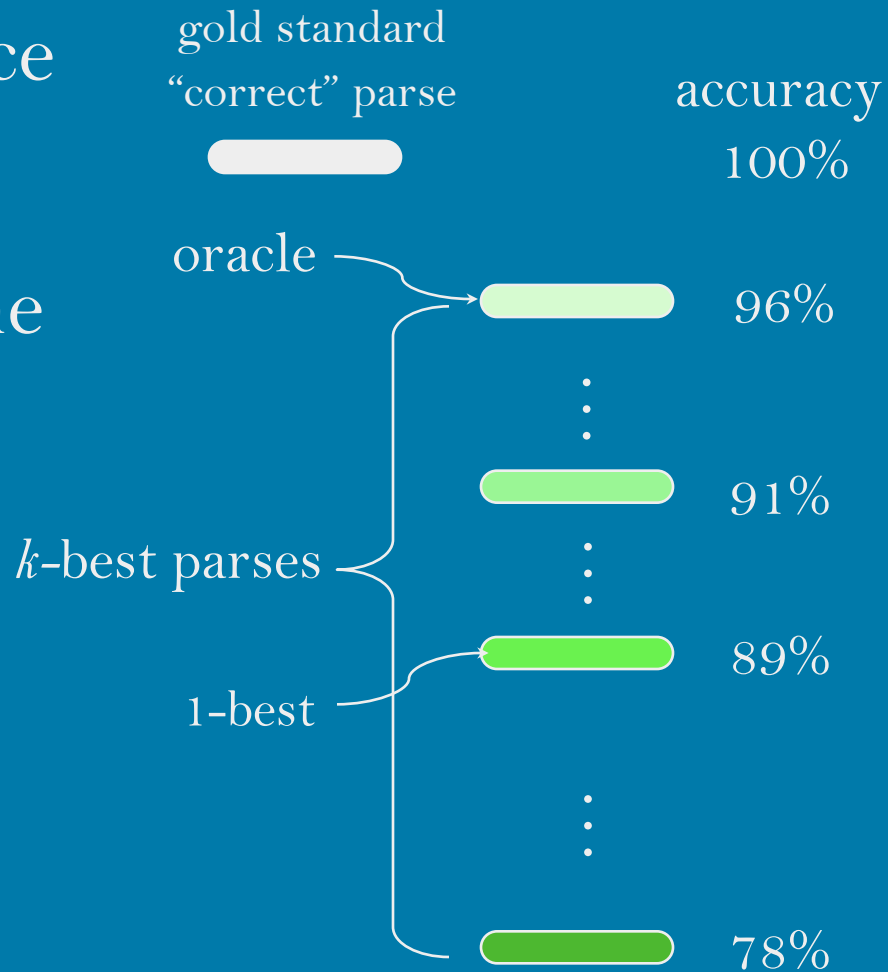
- **oracle reranking**: pick the best parse according to the gold-standard
- **real reranking**: pick the best parse according to the score function



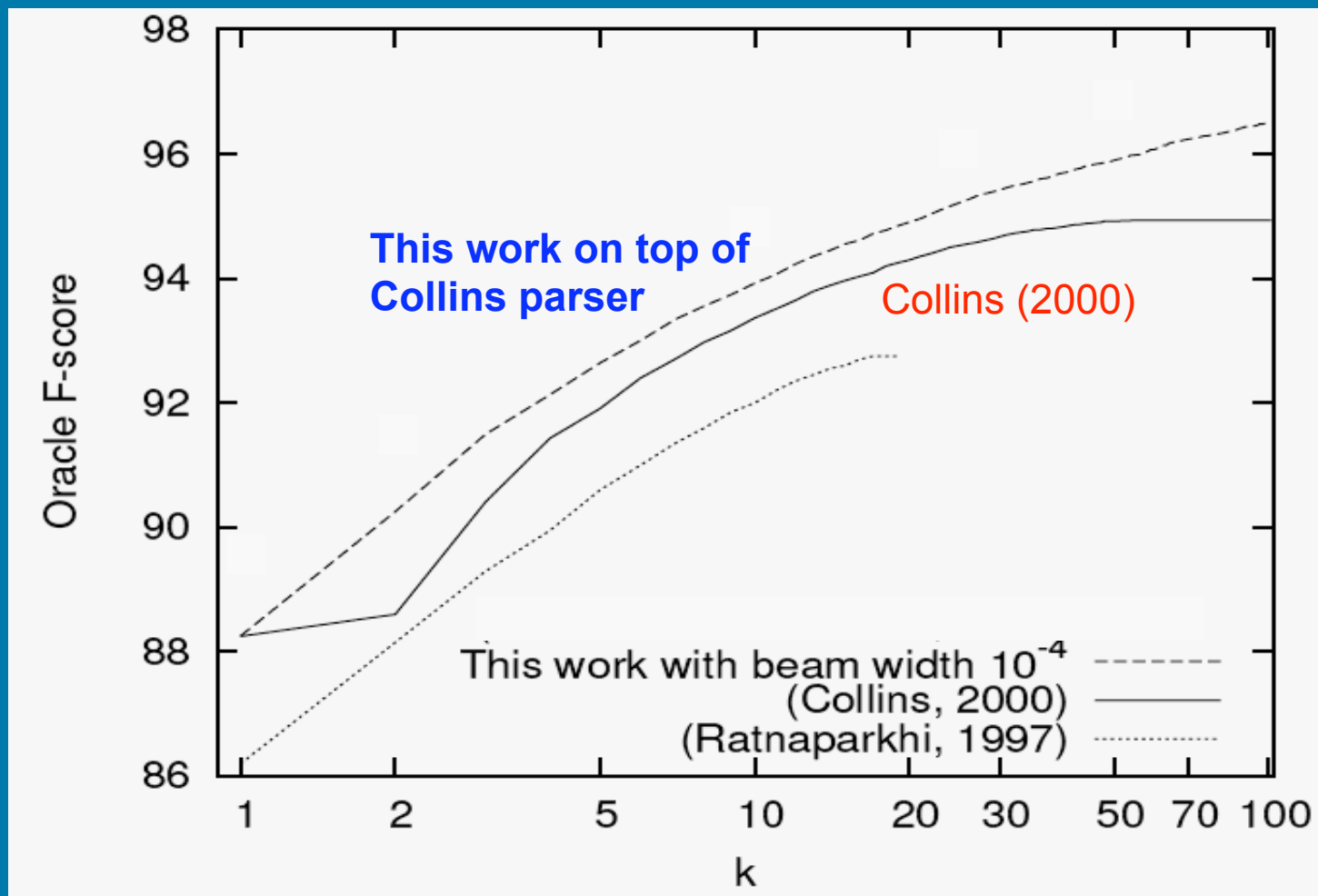
# Oracle Reranking

given  $k$  parses of a sentence

- **oracle reranking**: pick the best parse according to the gold-standard
- **real reranking**: pick the best parse according to the score function

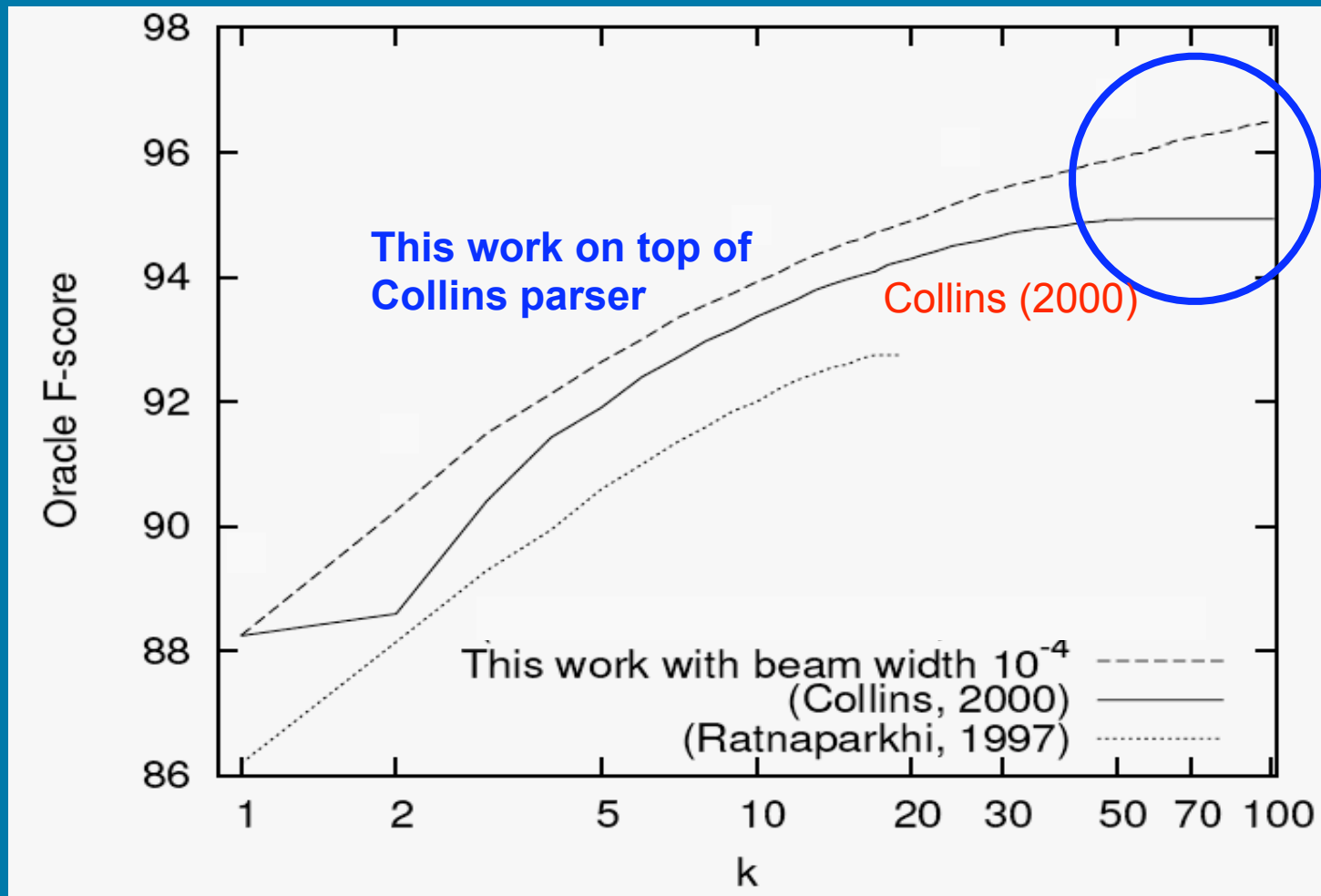


# Quality of the $k$ -best lists





# Quality of the $k$ -best lists



# Why are our $k$ -best lists better?

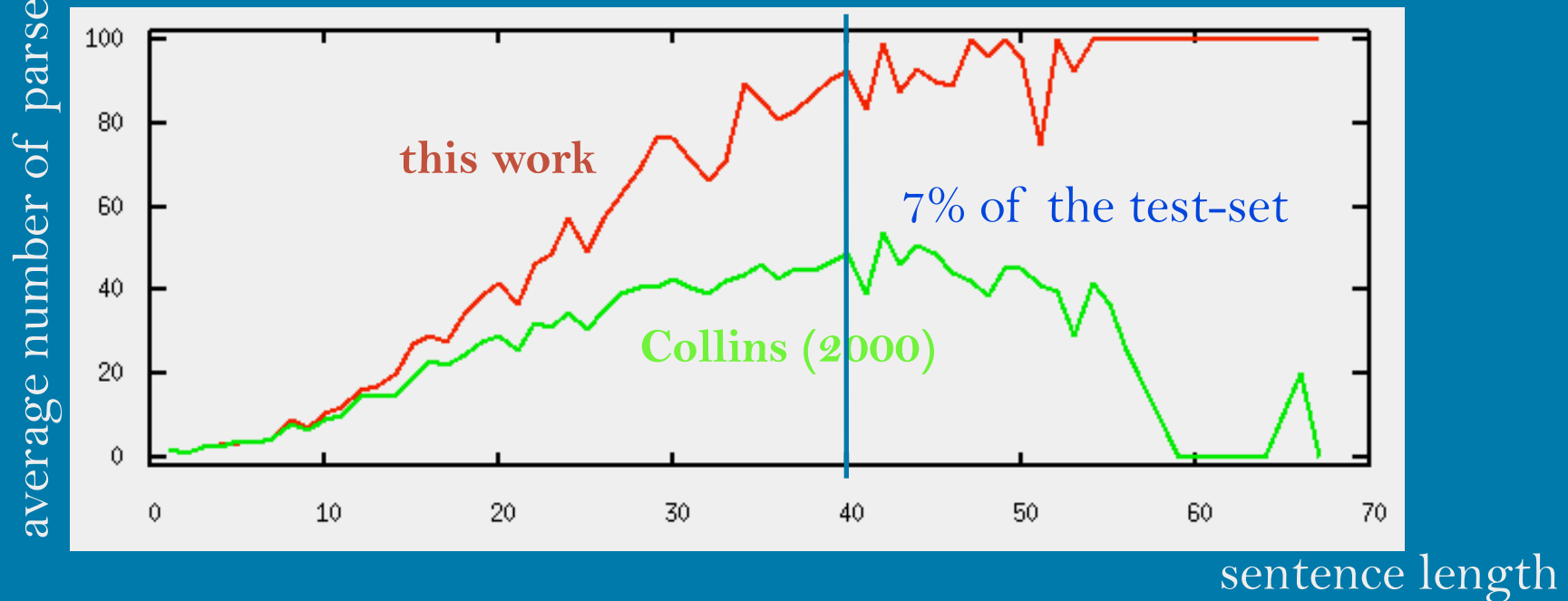
average number of parses for sentences of certain length



Collins (2000): turn down dynamic programming  
theoretically exponential time complexity;  
aggressive beam pruning to make it tractable in practice

# Why are our $k$ -best lists better?

average number of parses for sentences of certain length



Collins (2000): turn down dynamic programming  
theoretically exponential time complexity;  
aggressive beam pruning to make it tractable in practice

# Implemented in ...

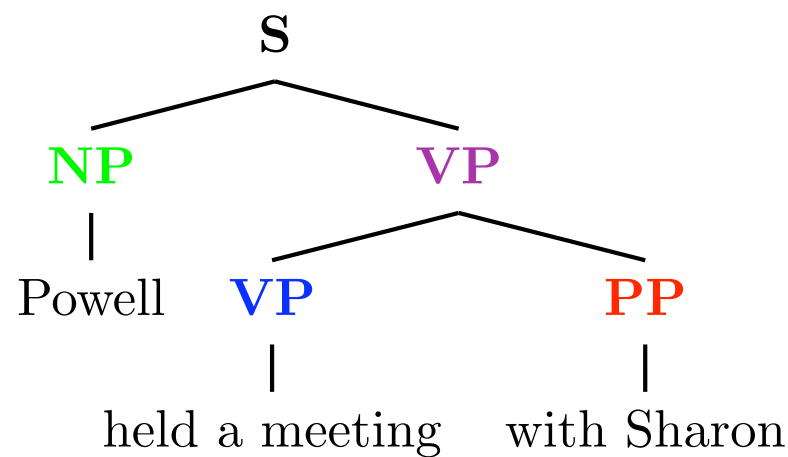
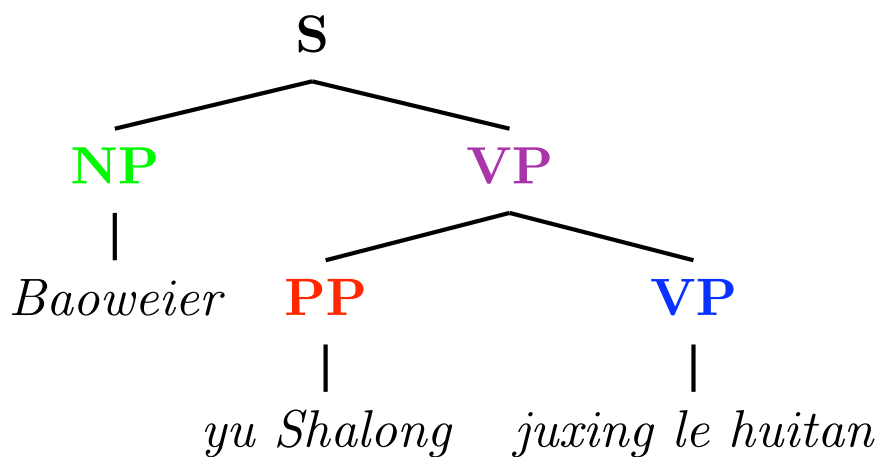
- state-of-the-art statistical parsers
  - Charniak parser (2005); Berkeley parser (2006)
  - McDonald et al. dependency parser (2005)
  - Microsoft Research (Redmond) dependency parser (2006)
- generic dynamic programming languages/packages
  - Dyna (Eisner et al., 2005) and Tiburon (May and Knight, 2006)
- state-of-the-art syntax-based translation systems
  - Hiero (Chiang, 2005)
  - ISI syntax-based system (2005)
  - CMU syntax-based system (2006)
  - BBN syntax-based system (2007)

# Applications in Machine Translation

# Syntax-based Translation

- synchronous context-free grammars (SCFGs)
- generating pairs of strings/trees simultaneously
- co-indexed nonterminal further rewritten as a unit

$S \rightarrow NP^{(1)} VP^{(2)}, \quad NP^{(1)} VP^{(2)}$   
 $VP \rightarrow PP^{(1)} VP^{(2)}, \quad VP^{(2)} PP^{(1)}$   
 $NP \rightarrow Baoweier, \quad Powell$   
...



# Translation as Parsing

- translation (“decoding”) => monolingual parsing
- parse the source input with the source projection
- build the corresponding target sub-strings in parallel

# Translation as Parsing

- translation (“decoding”) => monolingual parsing
- parse the source input with the source projection
- build the corresponding target sub-strings in parallel

S → NP<sup>(1)</sup> VP<sup>(2)</sup>, NP<sup>(1)</sup> VP<sup>(2)</sup>  
VP → PP<sup>(1)</sup> VP<sup>(2)</sup>, VP<sup>(2)</sup> PP<sup>(1)</sup>  
NP → *Baoweier*, Powell  
...

Baoweier

yu

Shalong

juxing le huitan



# Translation as Parsing

- translation (“decoding”) => monolingual parsing
- parse the source input with the source projection
- build the corresponding target sub-strings in parallel

S → NP<sup>(1)</sup> VP<sup>(2)</sup>,  
VP → PP<sup>(1)</sup> VP<sup>(2)</sup>,  
NP → *Baoweier*,  
...

Baoweier

yu

Shalong

juxing

le

huitan

# Translation as Parsing

- translation (“decoding”) => monolingual parsing
- parse the source input with the source projection
- build the corresponding target sub-strings in parallel

S → NP<sup>(1)</sup> VP<sup>(2)</sup>,  
VP → PP<sup>(1)</sup> VP<sup>(2)</sup>,  
NP → *Baoweier*,  
...

NP

Baoweier

PP

yu Shalong

VP

juxing le huitan

# Translation as Parsing

- translation (“decoding”) => monolingual parsing
- parse the source input with the source projection
- build the corresponding target sub-strings in parallel

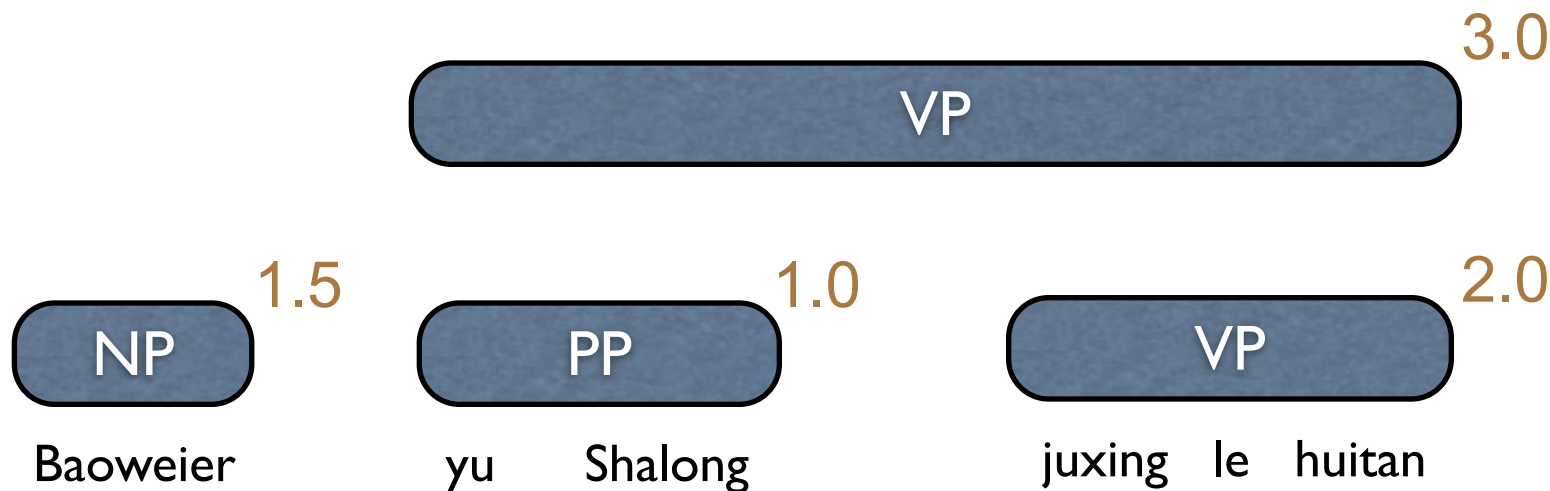
S → NP<sup>(1)</sup> VP<sup>(2)</sup>,  
VP → PP<sup>(1)</sup> VP<sup>(2)</sup>,  
NP → *Baoweier*,  
...



# Translation as Parsing

- translation (“decoding”) => monolingual parsing
- parse the source input with the source projection
- build the corresponding target sub-strings in parallel

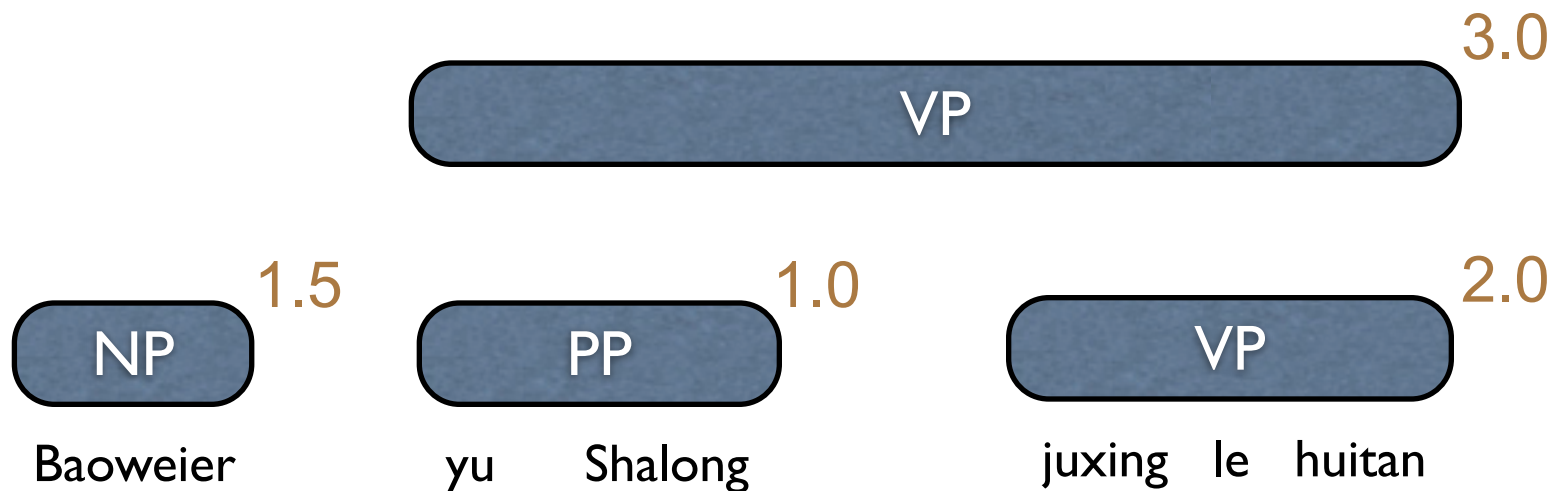
S → NP<sup>(1)</sup> VP<sup>(2)</sup>,  
VP → PP<sup>(1)</sup> VP<sup>(2)</sup>,  
NP → *Baoweier*,  
...



# Translation as Parsing

- translation (“decoding”) => monolingual parsing
- parse the source input with the source projection
- build the corresponding target sub-strings in parallel

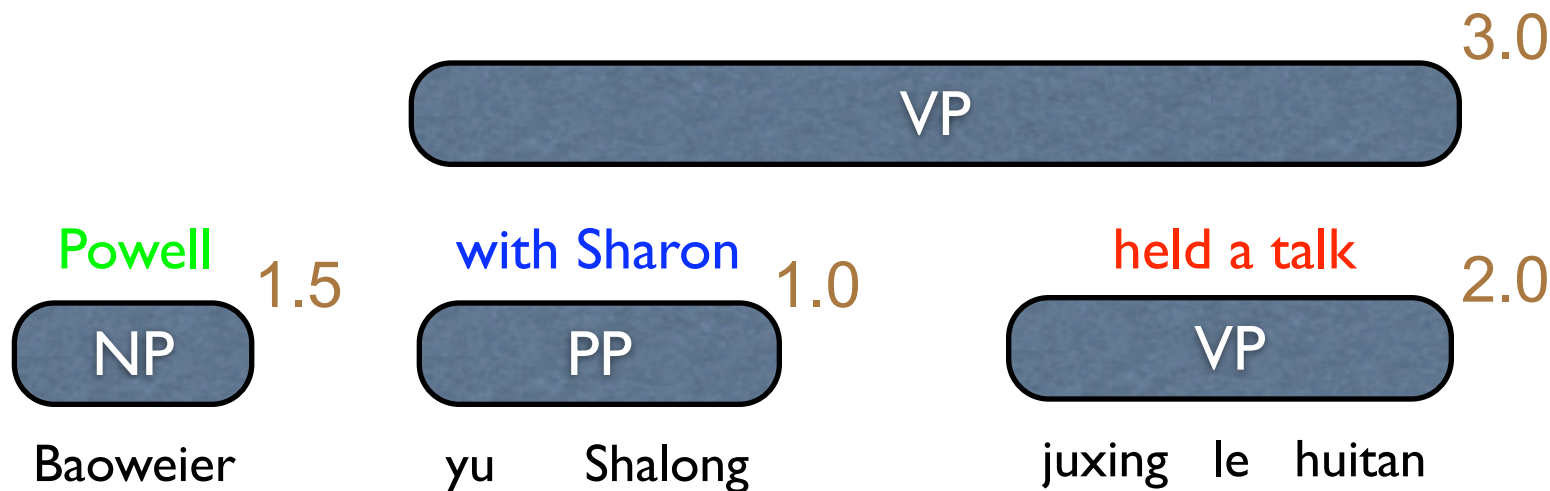
S → NP<sup>(1)</sup> VP<sup>(2)</sup>, NP<sup>(1)</sup> VP<sup>(2)</sup>  
VP → PP<sup>(1)</sup> VP<sup>(2)</sup>, VP<sup>(2)</sup> PP<sup>(1)</sup>  
NP → *Baoweier*, Powell  
...



# Translation as Parsing

- translation (“decoding”) => monolingual parsing
- parse the source input with the source projection
- build the corresponding target sub-strings in parallel

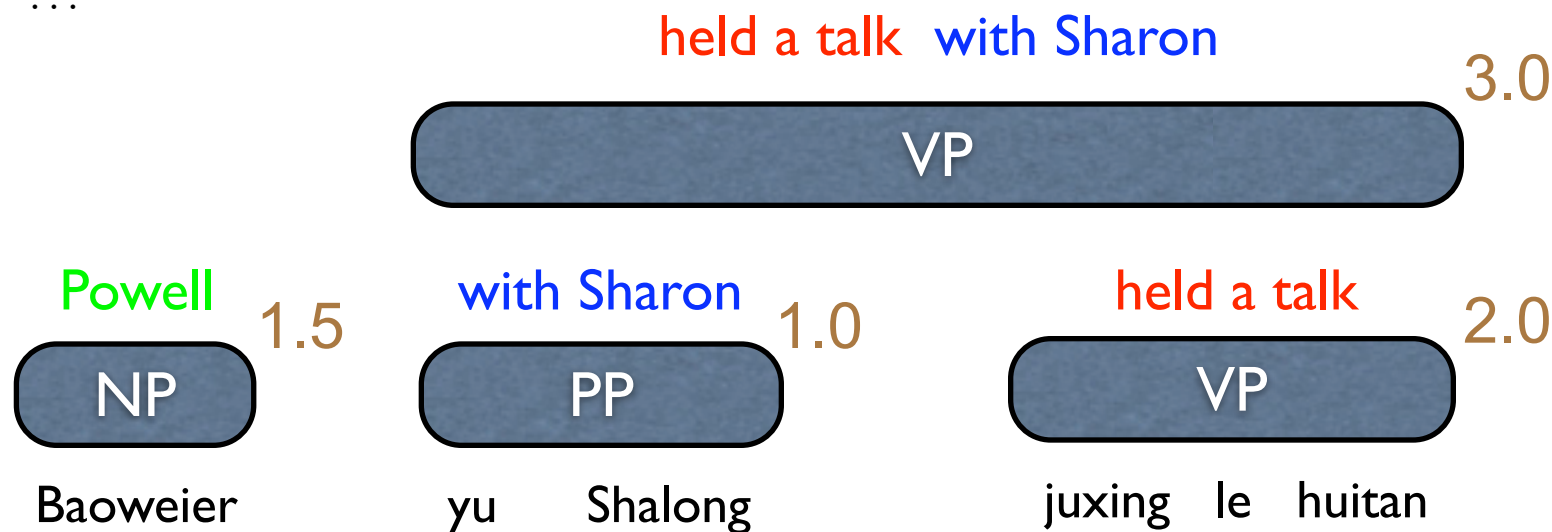
S → NP<sup>(1)</sup> VP<sup>(2)</sup>, NP<sup>(1)</sup> VP<sup>(2)</sup>  
VP → PP<sup>(1)</sup> VP<sup>(2)</sup>, VP<sup>(2)</sup> PP<sup>(1)</sup>  
NP → *Baoweier*, Powell  
...



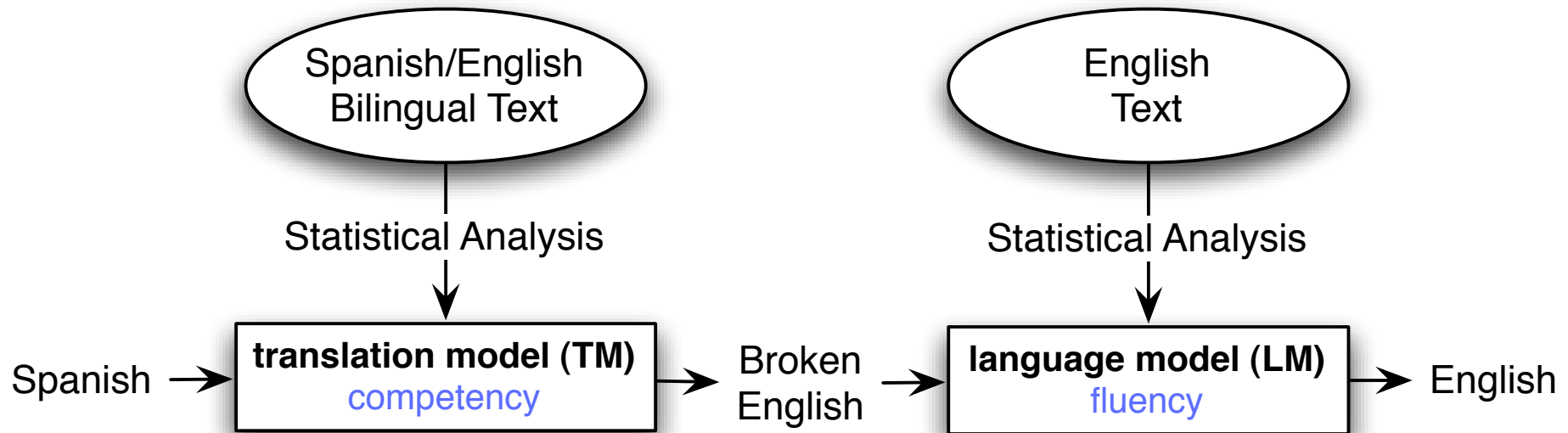
# Translation as Parsing

- translation (“decoding”) => monolingual parsing
- parse the source input with the source projection
- build the corresponding target sub-strings in parallel

S → NP<sup>(1)</sup> VP<sup>(2)</sup>, NP<sup>(1)</sup> VP<sup>(2)</sup>  
VP → PP<sup>(1)</sup> VP<sup>(2)</sup>, VP<sup>(2)</sup> PP<sup>(1)</sup>  
NP → *Baoweier*, Powell  
...



# Language Model: Rescoring

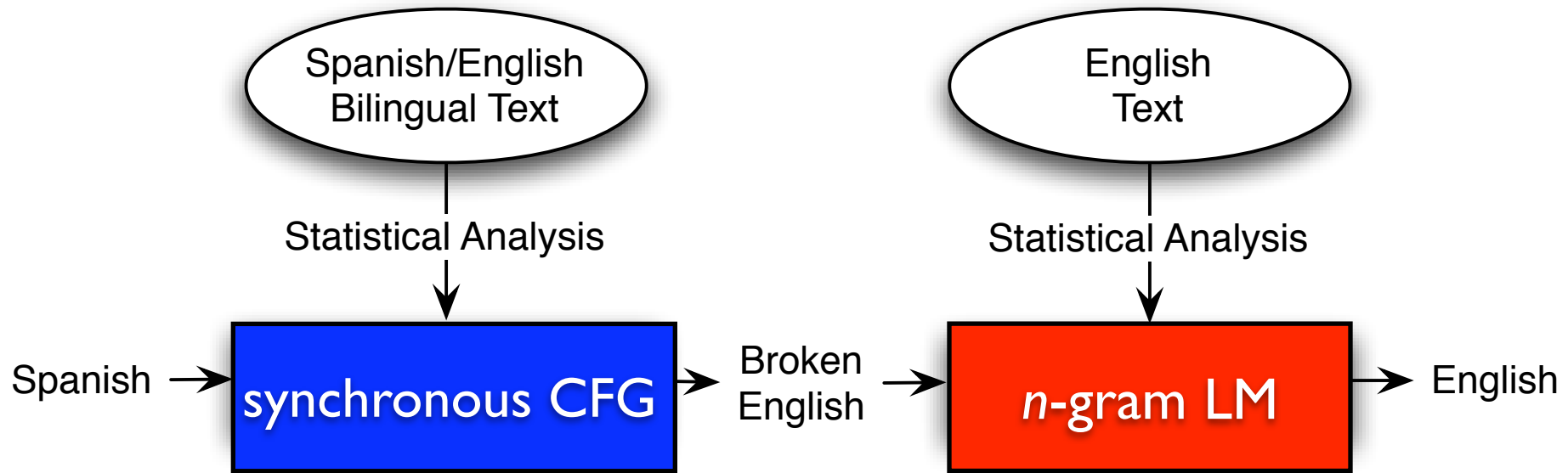


Que hambre tengo yo →  
What hunger have I  
Hungry I am so  
Have I that hunger  
I am so hungry  
How hunger have I  
→ I am so hungry

...



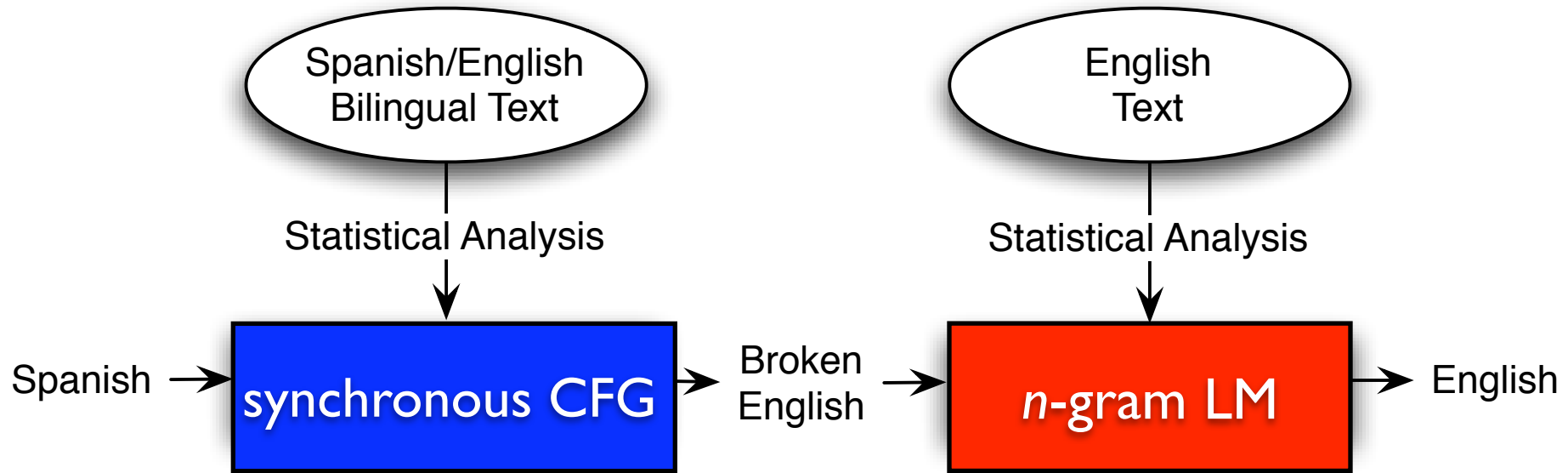
# Language Model: Rescoring



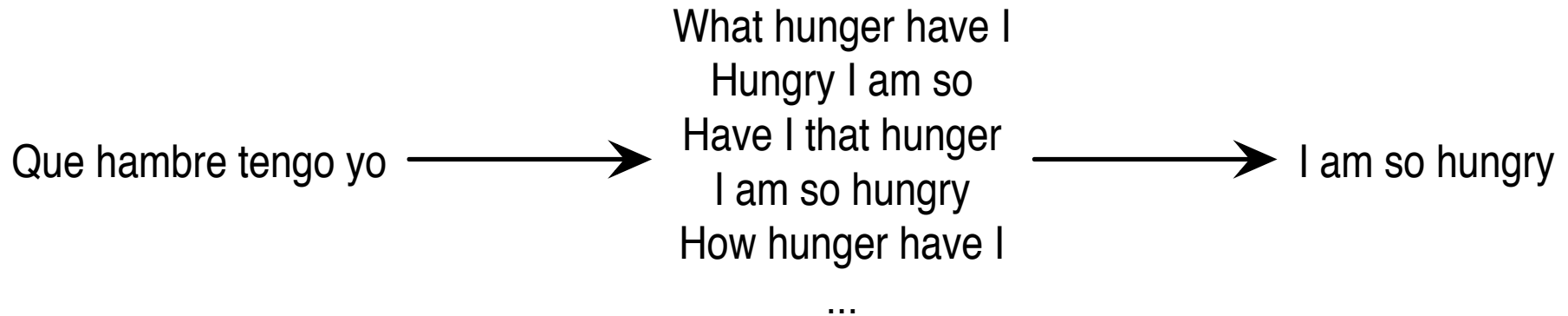
Que hambre tengo yo →  
What hunger have I  
Hungry I am so  
Have I that hunger  
I am so hungry  
How hunger have I  
→ I am so hungry

...

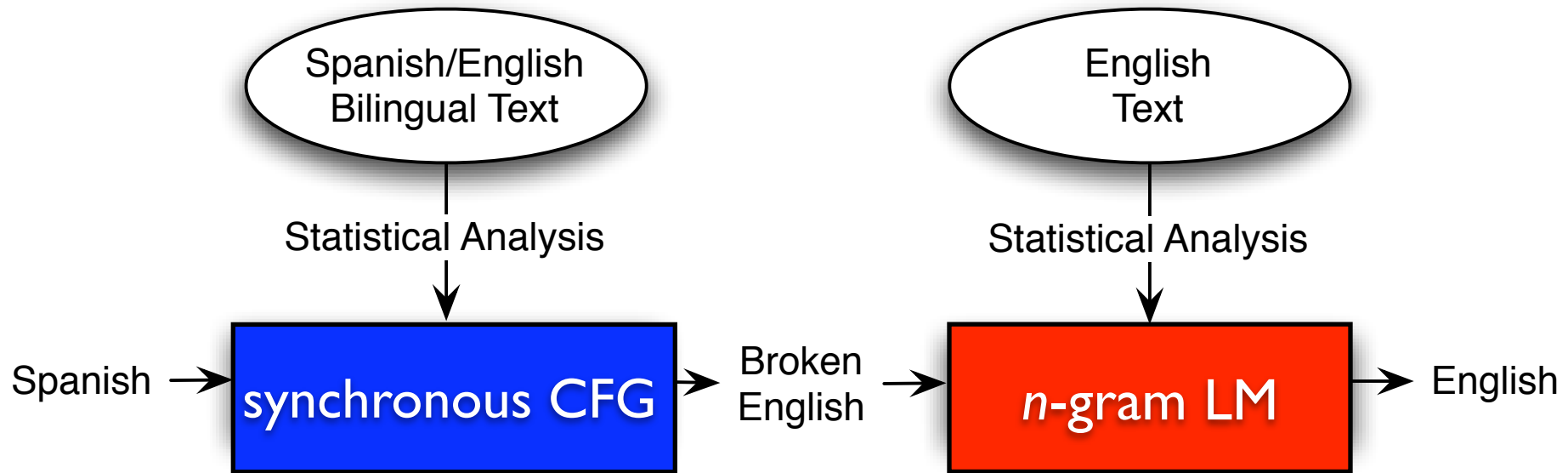
# Language Model: Rescoring



## *k*-best rescoring



# Language Model: Rescoring

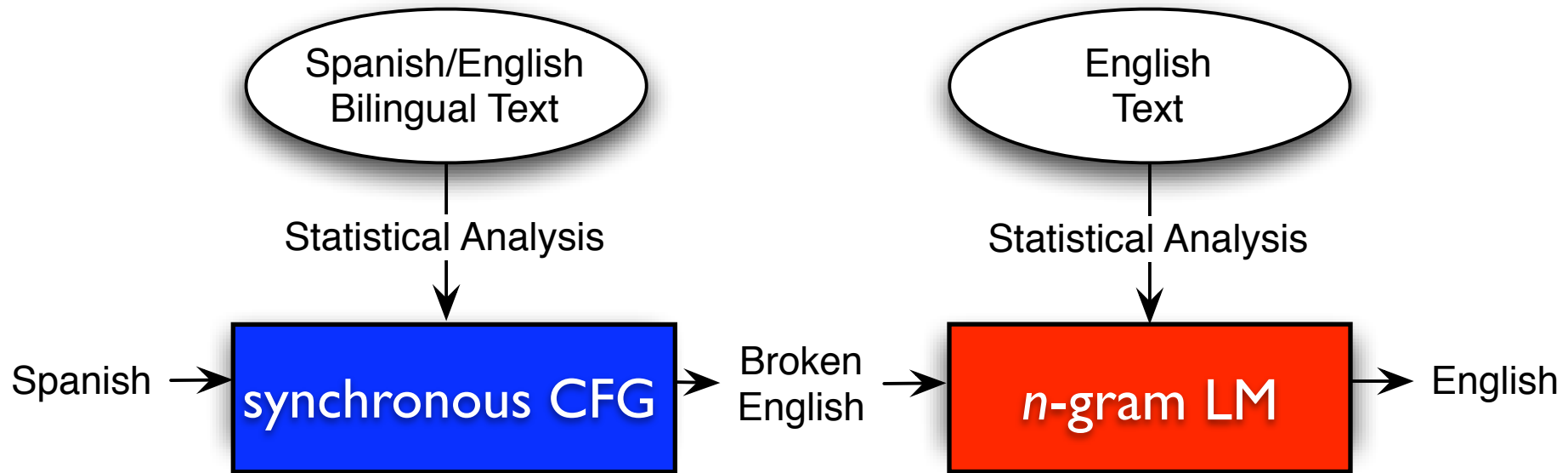


*k*-best rescoring

TM

	3.7	What hunger have I	
	3.8	Hungry I am so	
Que hambre tengo yo	4.1	Have I that hunger	→ I am so hungry
	4.5	I am so hungry	
	7.2	How hunger have I	
		...	

# Language Model: Rescoring



*k*-best rescoring

TM

LM

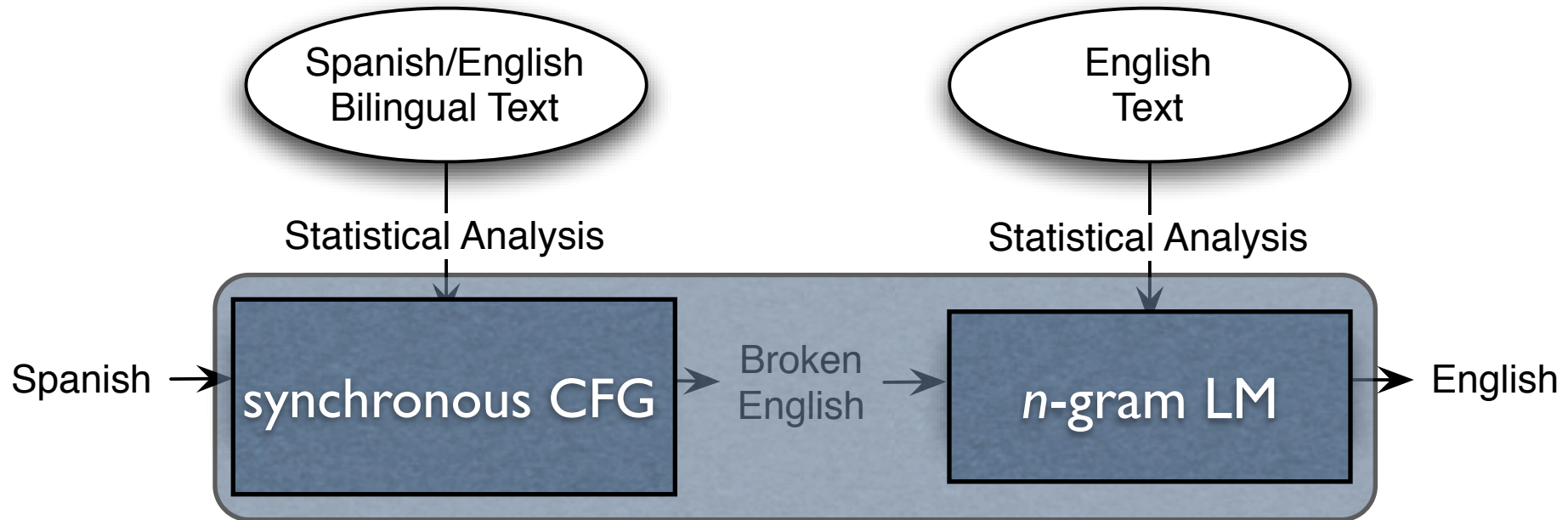
	3.7	What hunger have I	5.4	
	3.8	Hungry I am so	7.0	
Que hambre tengo yo →	4.1	Have I that hunger	9.8	→ I am so hungry
	4.5	I am so hungry	0.2	
	7.2	How hunger have I	8.7	

...

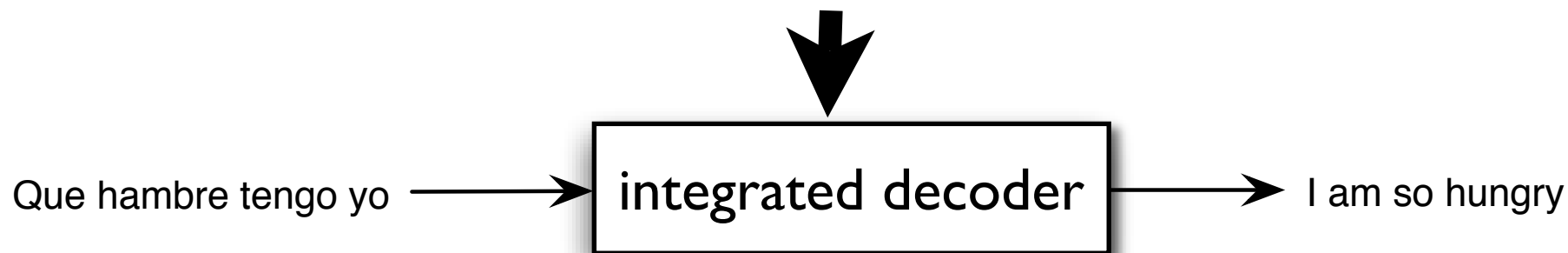
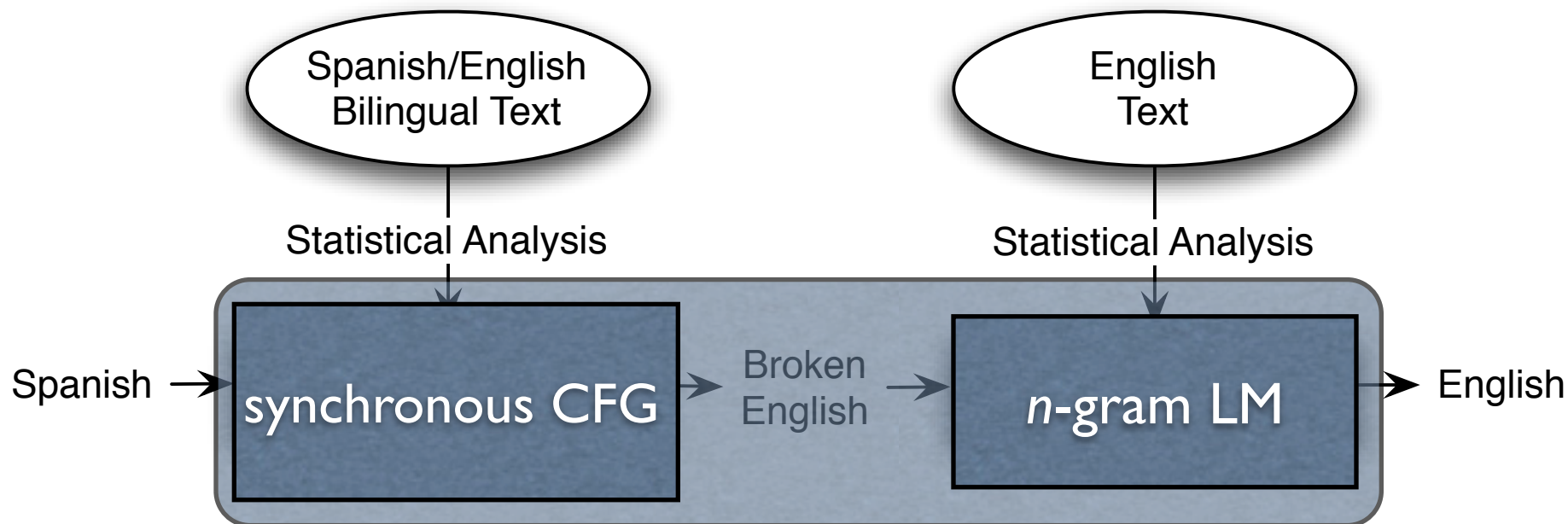
# k-best rescoring results

- The ISI syntax-based translation system
  - currently the best performing system on Chinese to English task in NIST evaluations
  - based on synchronous grammars
  - translation model (TM) only: BLEU score 24.45
  - rescoring with trigram LM on 25000-best list: 34.58

# Language Model: Integration



# Language Model: Integration



computationally challenging! ☹️

# Integrated Decoding Results

- The ISI syntax-based translation system
  - currently the best performing system on Chinese to English task in NIST evaluations
  - based on synchronous grammars
  - translation model (TM) only: BLEU score 24.45
  - rescoring with trigram LM on 25000-best list: 34.58
  - trigram integrated decoding: 38.44

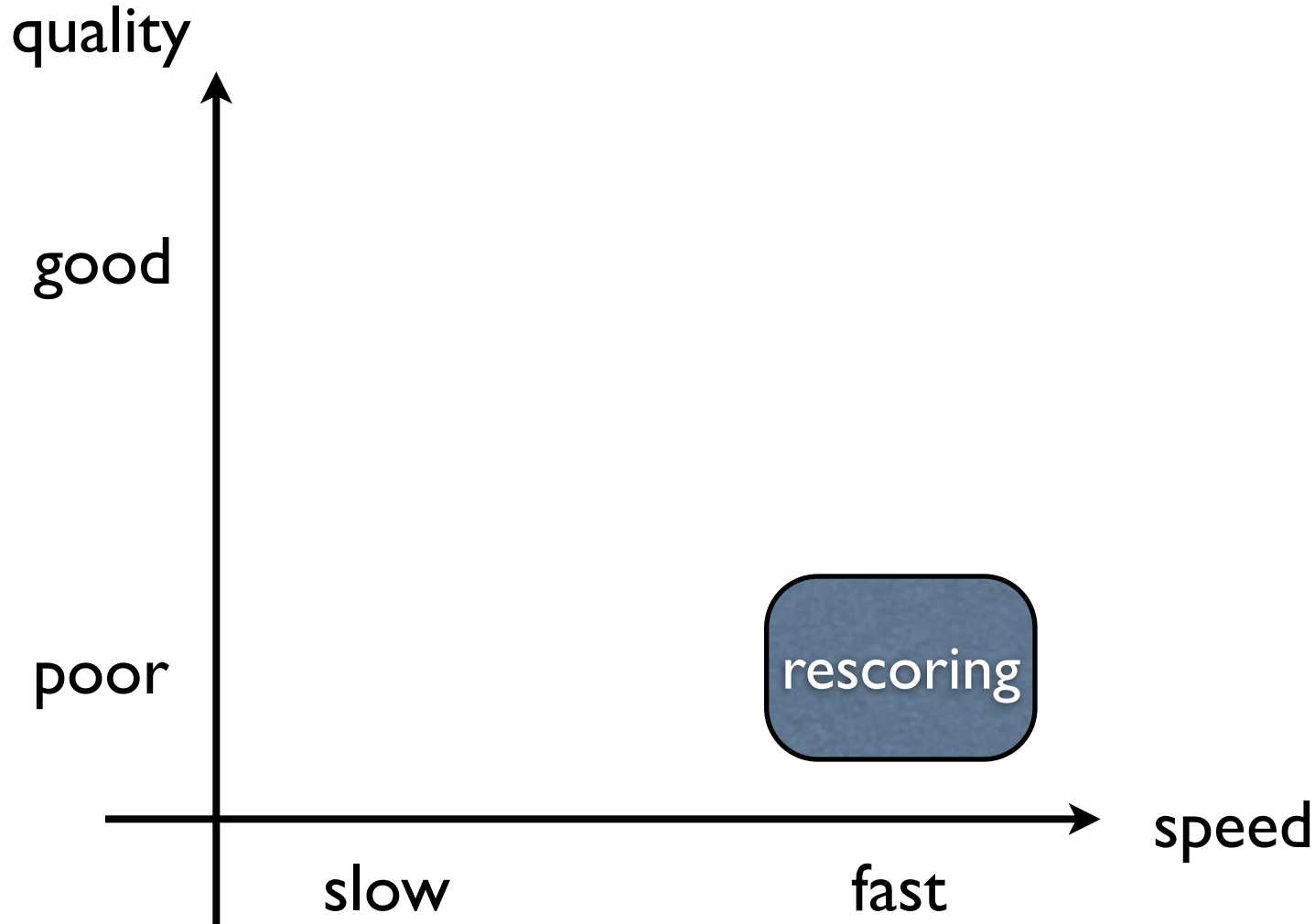


# Integrated Decoding Results

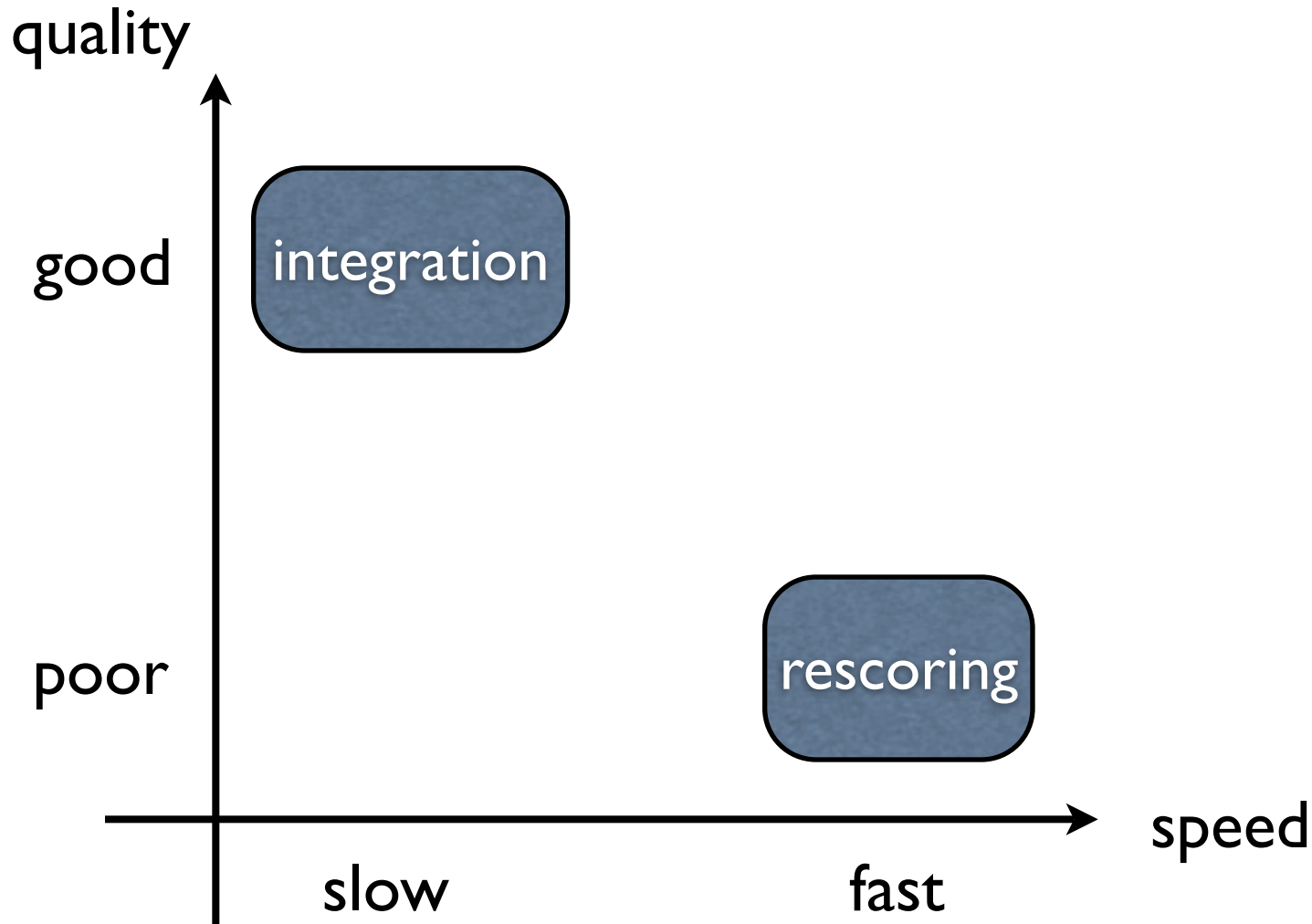
- The ISI syntax-based translation system
  - currently the best performing system on Chinese to English task in NIST evaluations
  - based on synchronous grammars
  - translation model (TM) only: BLEU score 24.45
  - rescoring with trigram LM on 25000-best list: 34.58
  - trigram integrated decoding: 38.44

**but very slow!**

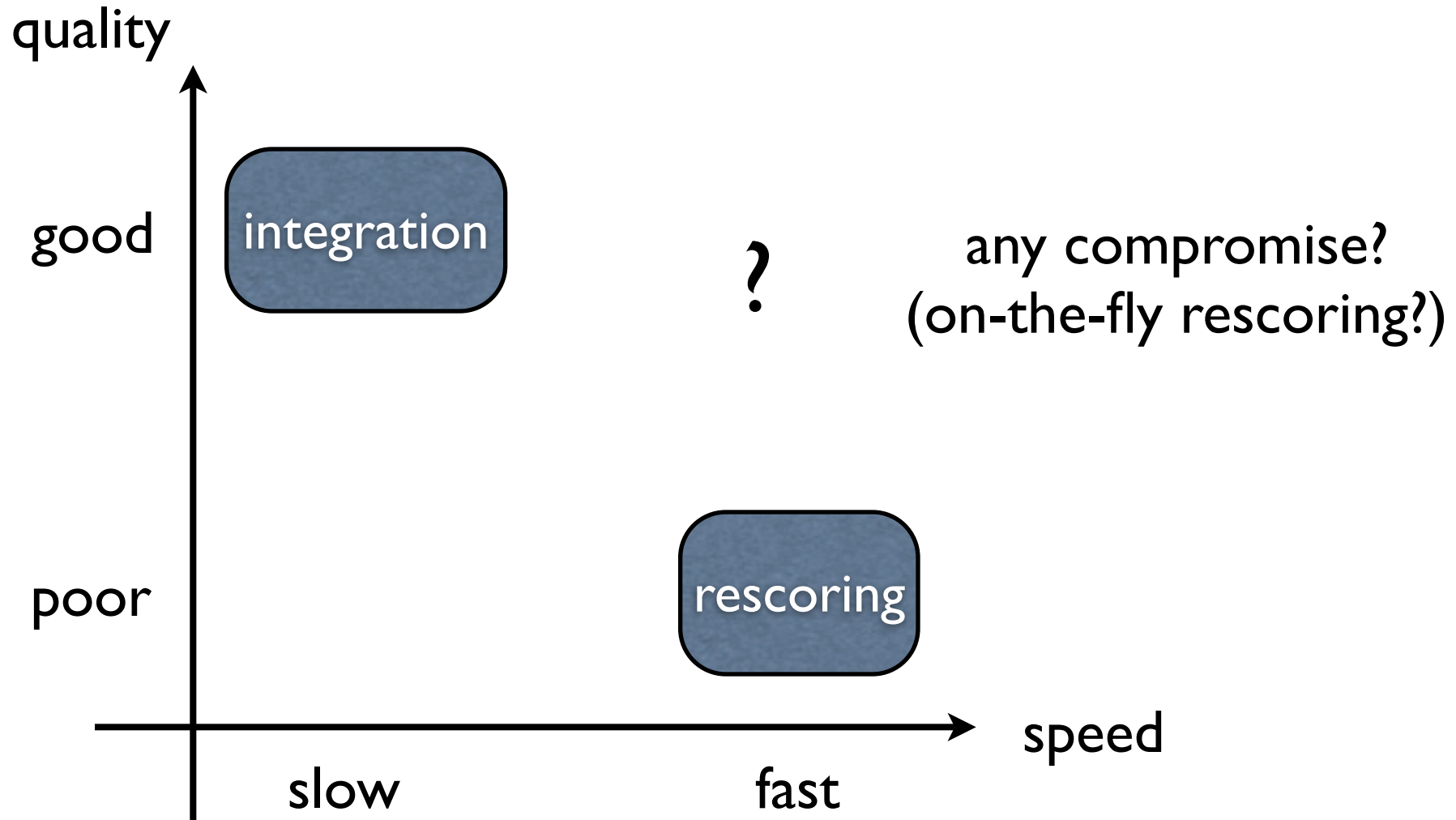
# Rescoring vs. Integration



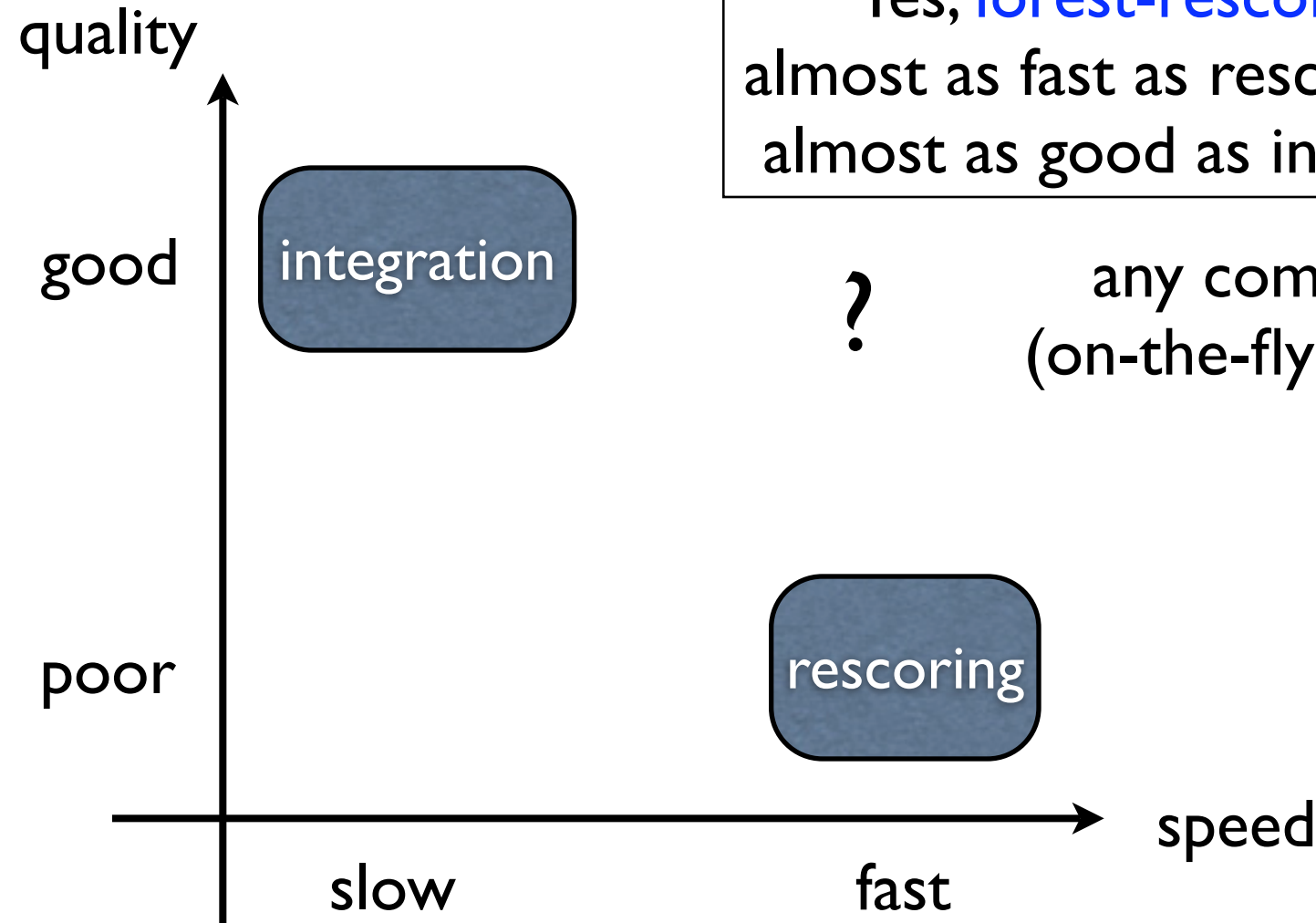
# Rescoring vs. Integration



# Rescoring vs. Integration



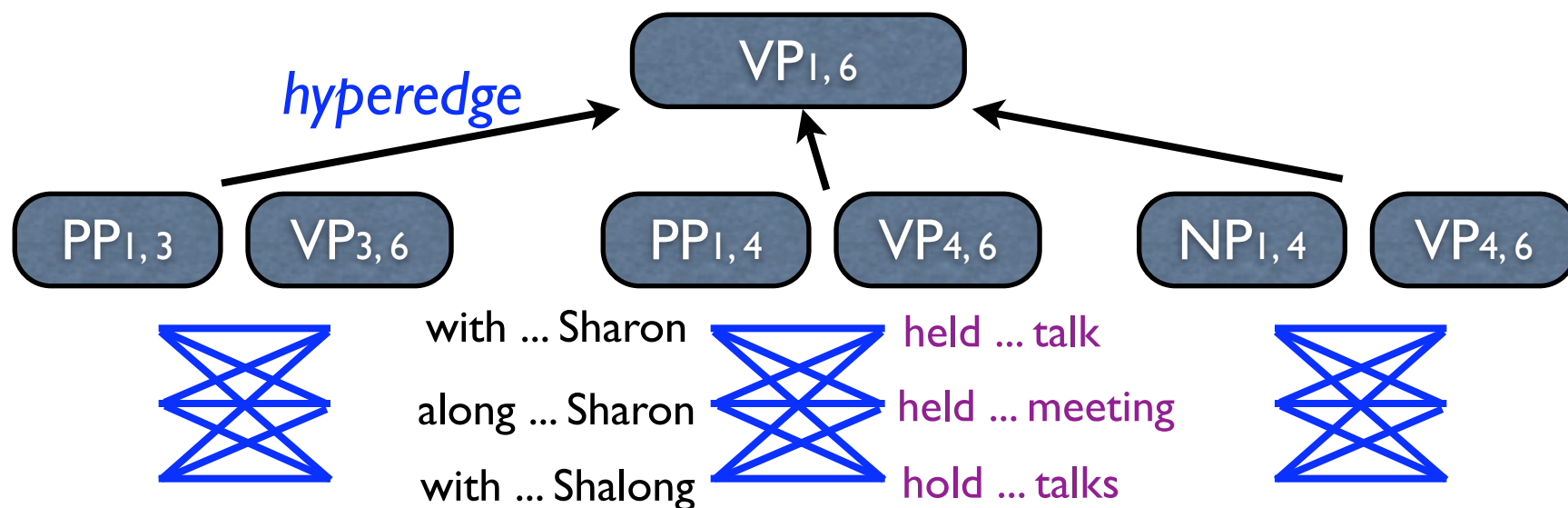
# Rescoring vs. Integration



Yes, **forest-rescoring** -- almost as fast as rescoring, and almost as good as integration

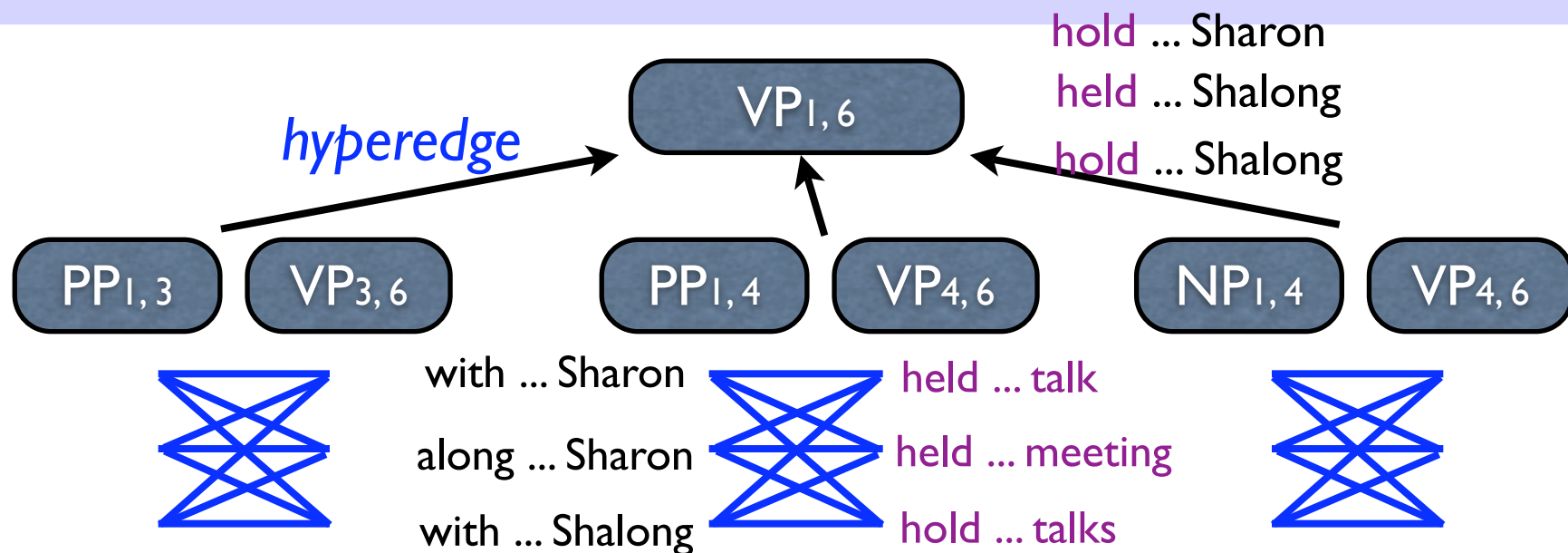
? any compromise?  
(on-the-fly rescoring?)

# Why Integration is Slow?



- split each node into +LM items (w/ boundary words)
- beam search: only keep **top-k** +LM items at each node
- but there are many ways to derive each node
- can we avoid enumerating all combinations?

# Why Integration is Slow?

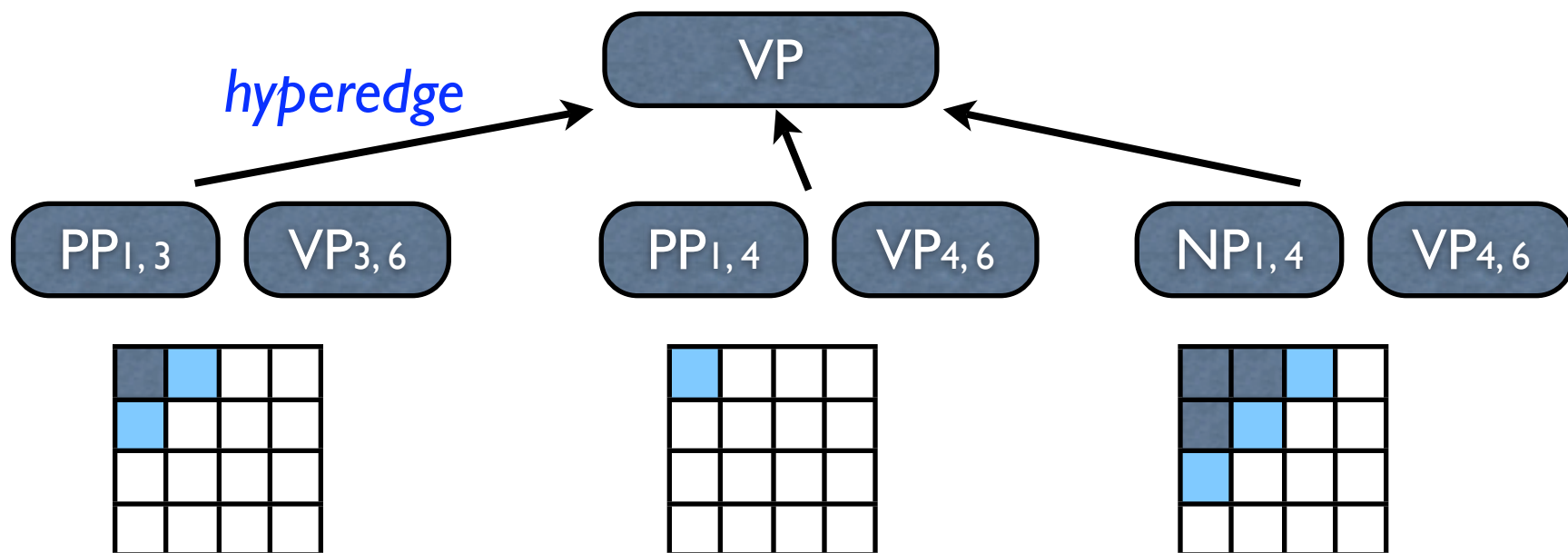


- split each node into +LM items (w/ boundary words)
- beam search: only keep **top-k** +LM items at each node
- but there are many ways to derive each node
- can we avoid enumerating all combinations?

# Forest Rescoring

$k$ -best parsing  
Algorithm 2

with LM cost, we can only  
do  $k$ -best **approximately**.



process all hyperedges **simultaneously!**  
significant savings of computation



# Forest Rescoring Results

- on the Hiero system (Chiang, 2005)
  - ~10 fold speed-up at the same level of BLEU
- on my syntax-directed system (Huang et al., 2006)
  - ~10 fold speed-up at the same level of search-error
- on a typical phrase-based system (Pharaoh)
  - ~30 fold speed-up at the same level of search-error
  - ~100 fold speed-up at the same level of BLEU
- also used in [see my ACL '07 paper for details](#)
  - ISI, CMU, and BBN syntax-based systems

# Conclusions

- monotonic hypergraph formulation
  - the  $k$ -best **derivations** problem
- $k$ -best Algorithms
  - Algorithm 0 (naïve) to Algorithm 3 (lazy)
- experimental results
  - efficiency
  - accuracy (effectively searching over larger space)
- applications in machine translation
  - $k$ -best rescoring and forest rescoring

**Thank you!**  
**谢谢!**

**Questions?**  
**Comments?**

Thank you!  
谢谢!



Questions?  
Comments?

# Thank you!

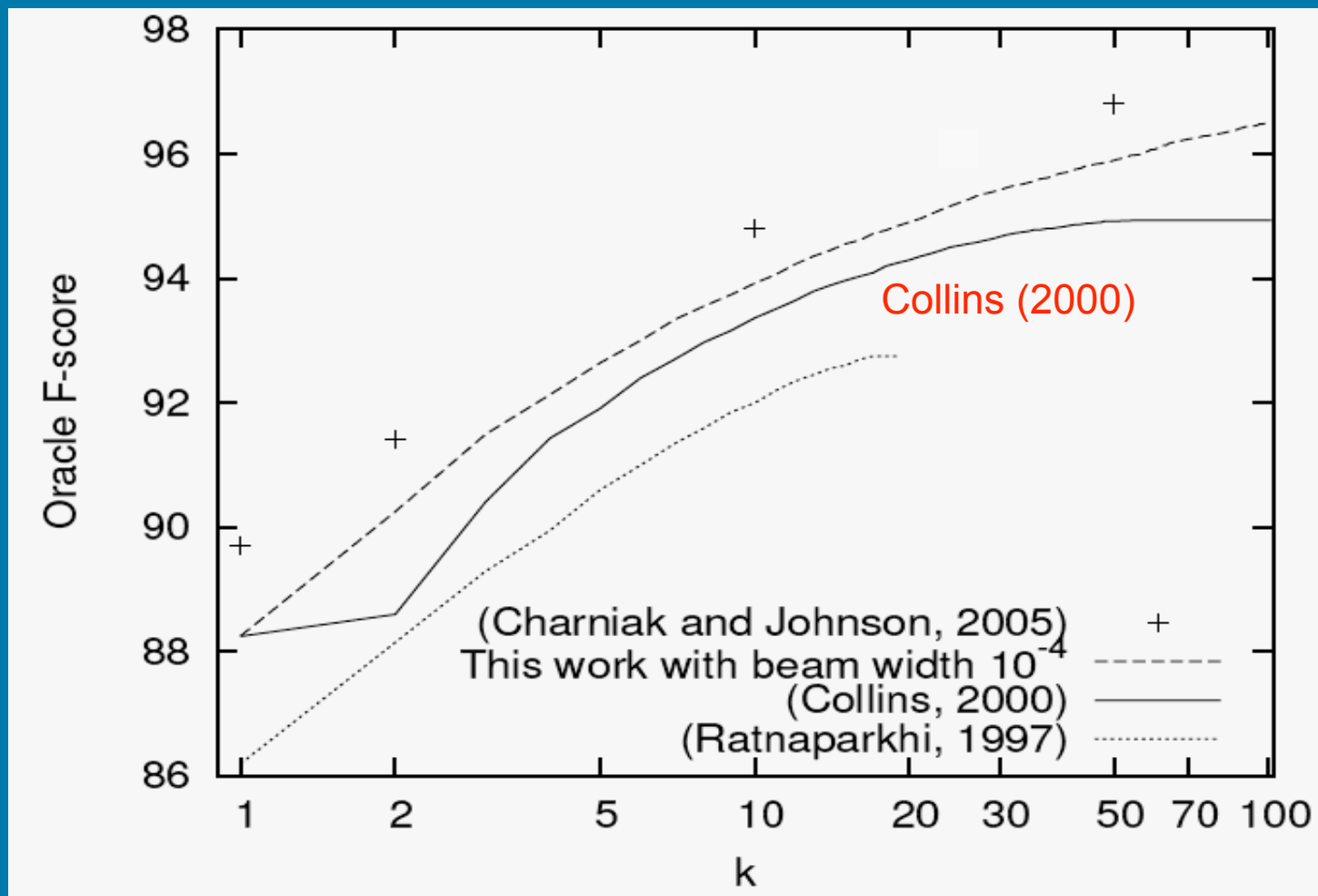
# 谢谢!



Questions?  
Comments?

- Liang Huang and David Chiang (2005).  
Better k-best Parsing.  
In Proceedings of IWPT, Vancouver, B.C.
- Liang Huang and David Chiang (2007).  
Forest Rescoring: Faster Decoding with Integrated Language Models.  
In Proceedings of ACL, Prague, Czech Rep.

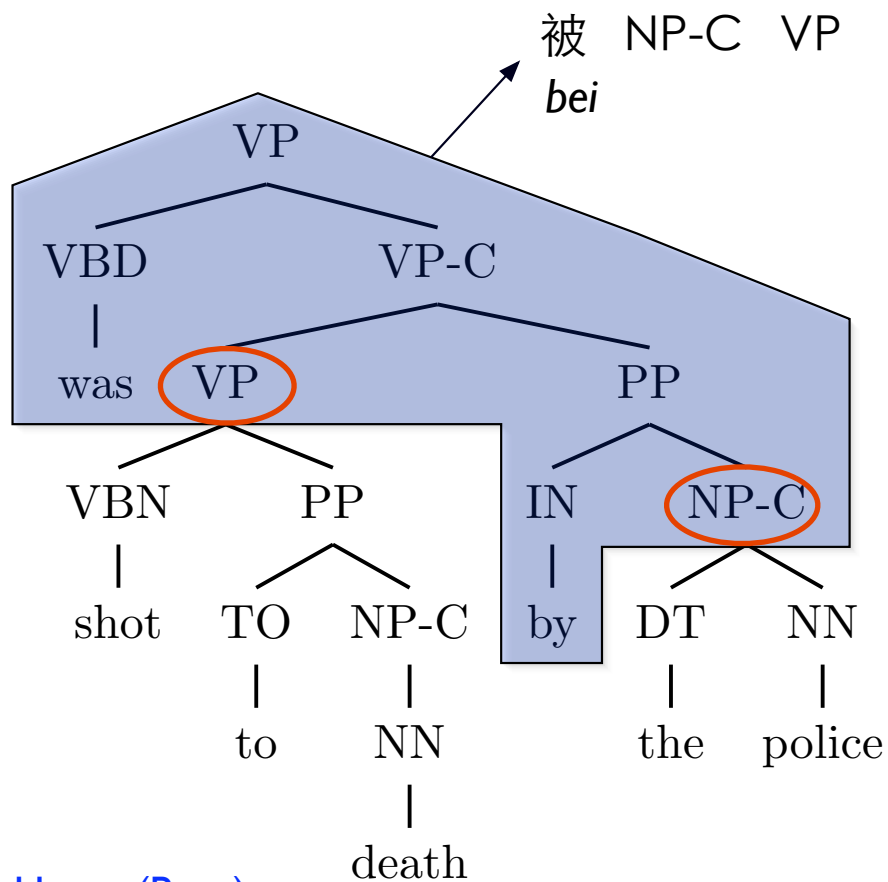
# Quality of the $k$ -best lists



# Syntax-based Experiments

# Tree-to-String System

- syntax-directed, English to Chinese (Huang, Knight, Joshi, 2006)
- the reverse direction is found in (Liu et al., 2006)



synchronous tree-  
substitution grammars (STSG)

(Galley et al., 2004; Eisner, 2003)

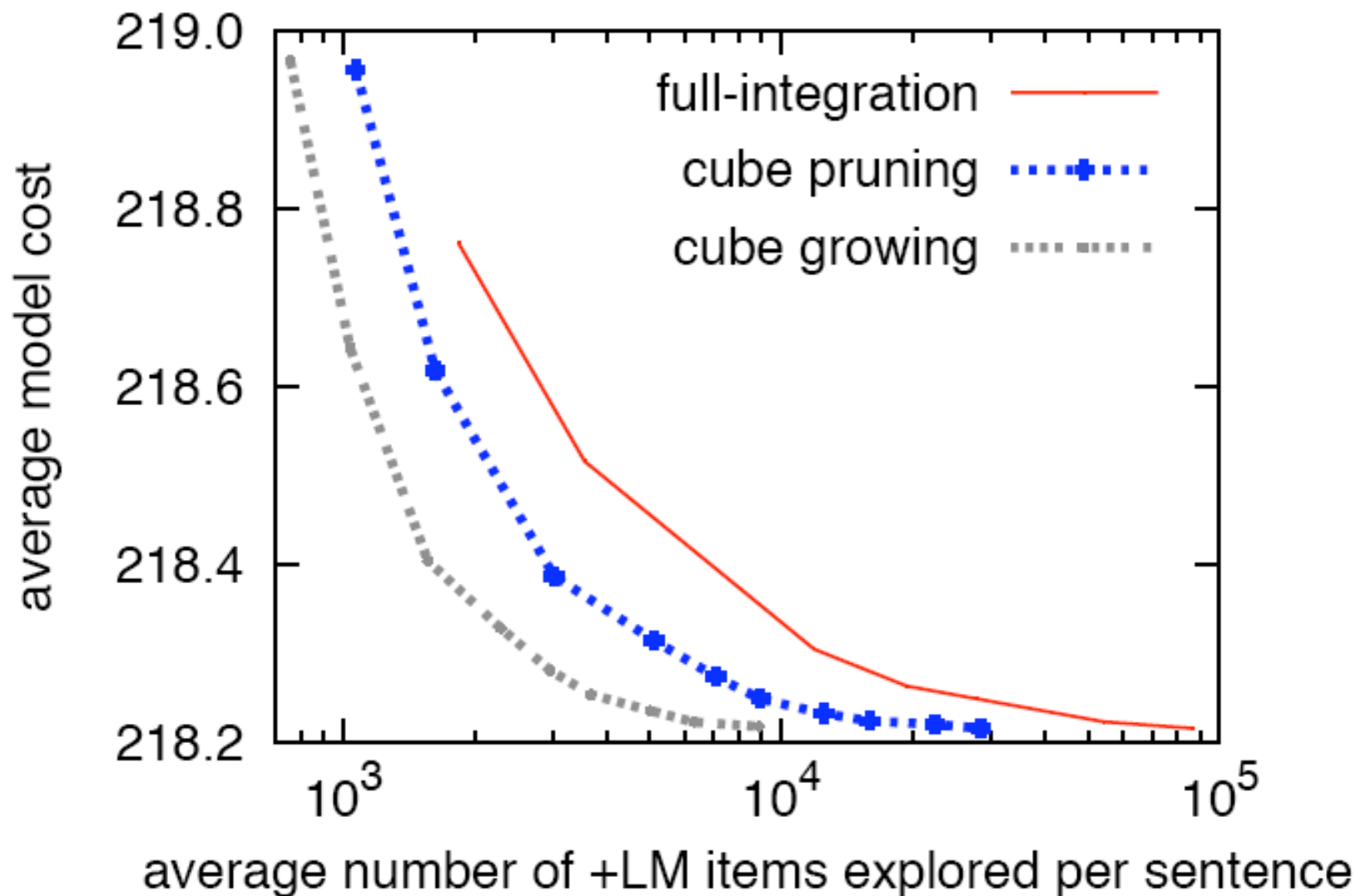
related to

STAG (Shieber/Schabes, 90)

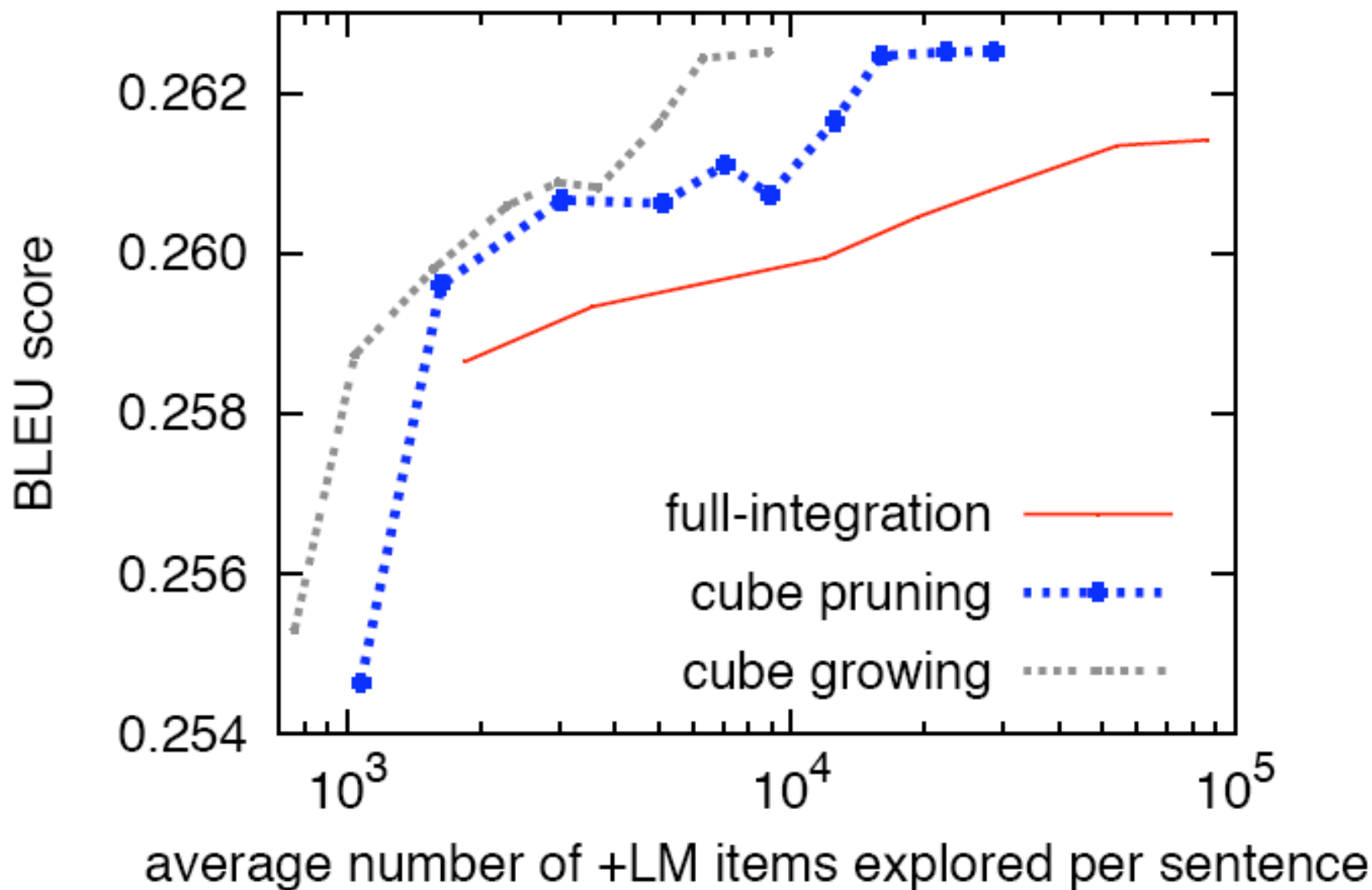
tested on 140 sentences  
slightly better BLEU scores  
than Pharaoh



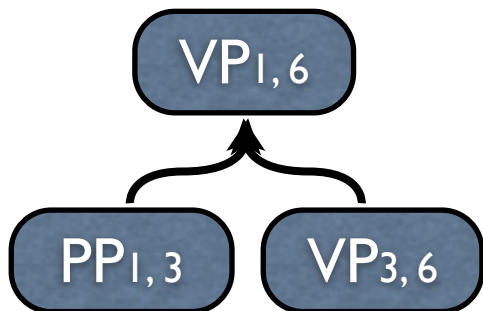
# Speed vs. Search Quality



# Speed vs. Translation Accuracy



# Cube Pruning



(PP with \* Sharon)  
 (PP along \* Sharon)  
 (PP with \* Shalong)

monotonic grid?

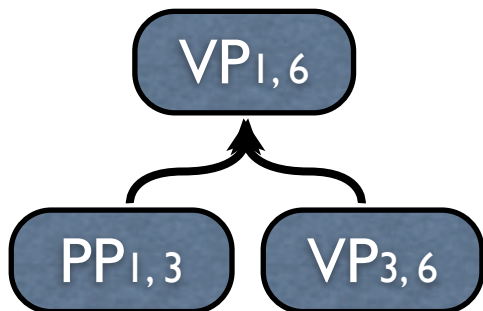
(VP held \* meeting)  
<sub>3,6</sub>

(VP held \* talk)  
<sub>3,6</sub>

(VP hold \* conference)  
<sub>3,6</sub>

	1.0	3.0	8.0
1.0	2.0	4.0	9.0
1.1	2.1	4.1	9.1
3.5	4.5	6.5	11.5

# Cube Pruning



non-monotonic grid  
due to LM combo costs

(VP<sub>3,6</sub> held \* meeting)

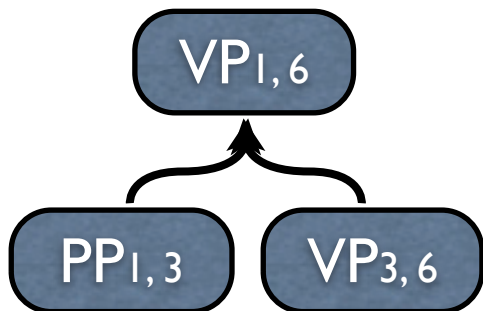
(VP<sub>3,6</sub> held \* talk)

(VP<sub>3,6</sub> hold \* conference)

(PP<sub>1,3</sub> with \* Sharon)  
(PP<sub>1,3</sub> along \* Sharon)  
(PP<sub>1,3</sub> with \* Shalong)

	1.0	3.0	8.0
1.0	2.0 + 0.5	4.0 + 5.0	9.0 + 0.5
1.1	2.1 + 0.3	4.1 + 5.4	9.1 + 0.3
3.5	4.5 + 0.6	6.5 + 10.5	11.5 + 0.6

# Cube Pruning



bigram (meeting, with)

(PP<sub>1,3</sub> with \* Sharon)

(PP<sub>1,3</sub> along \* Sharon)

(PP<sub>1,3</sub> with \* Shalong)

non-monotonic grid  
due to LM combo costs

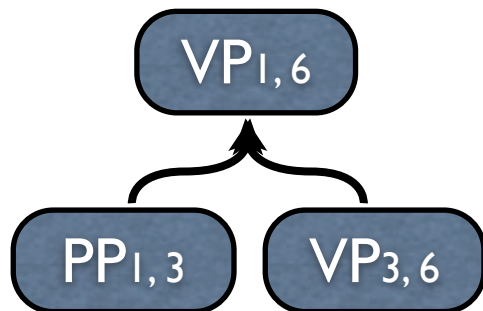
(VP<sub>3,6</sub> held \* meeting)

(VP<sub>3,6</sub> held \* talk)

(VP<sub>3,6</sub> hold \* conference)

	1.0	3.0	8.0
1.0	2.0 + 0.5	4.0 + 5.0	9.0 + 0.5
1.1	2.1 + 0.3	4.1 + 5.4	9.1 + 0.3
3.5	4.5 + 0.6	6.5 + 10.5	11.5 + 0.6

# Cube Pruning



non-monotonic grid  
due to LM combo costs

(VP<sub>3,6</sub> held \* meeting)

(VP<sub>3,6</sub> held \* talk)

(VP<sub>3,6</sub> hold \* conference)

(PP<sub>1,3</sub> with \* Sharon)

(PP<sub>1,3</sub> along \* Sharon)

(PP<sub>1,3</sub> with \* Shalong)

	1.0	3.0	8.0
1.0	2.5	9.0	9.5
1.1	2.4	9.5	9.4
3.5	5.1	17.0	12.1

# Cube Pruning

*k*-best parsing  
Algorithm 1

- a priority queue of candidates
- extract the best candidate

(PP with \* Sharon)  
1,3

(PP along \* Sharon)  
1,3

(PP with \* Shalong)  
1,3

(VP held \* meeting)  
3,6

(VP held \* talk)  
3,6

(VP hold \* conference)  
3,6

	1.0	3.0	8.0	
(VP held * meeting) 3,6	1.0	2.5	9.0	9.5
(VP held * talk) 3,6	1.1	2.4	9.5	9.4
(VP hold * conference) 3,6	3.5	5.1	17.0	12.1

# Cube Pruning

## k-best parsing Algorithm I

- a priority queue of candidates
- extract the best candidate
- push the two successors

(PP with \* Sharon)  
1,3

(PP along \* Sharon)  
1,3

(PP with \* Shalong)  
1,3

(VP held \* meeting)  
3,6

(VP held \* talk)  
3,6

(VP hold \* conference)  
3,6

	1.0	3.0	8.0
1.0	2.5	9.0	9.5
1.1	2.4	9.5	9.4
3.5	5.1	17.0	12.1



# Cube Pruning

## k-best parsing Algorithm I

- a priority queue of candidates
- extract the best candidate
- push the two successors

(PP with \* Sharon)  
1,3

(PP along \* Sharon)  
1,3

(PP with \* Shalong)  
1,3

(VP held \* meeting)  
3,6

(VP held \* talk)  
3,6

(VP hold \* conference)  
3,6

	1.0	3.0	8.0	
(VP held * meeting) 3,6	1.0	2.5	9.0	9.5
(VP held * talk) 3,6	1.1	2.4	9.5	9.4
(VP hold * conference) 3,6	3.5	5.1	17.0	12.1

# Cube Pruning

items are popped out-of-order

**solution:** keep a buffer of pop-ups

2.5 2.4 5.1

(PP with \* Sharon)  
1,3

(PP along \* Sharon)  
1,3

(PP with \* Sharon)  
1,3

(VP held \* meeting)  
3,6

(VP held \* talk)  
3,6

(VP hold \* conference)  
3,6

	1.0	3.0	8.0
(VP held * meeting) 3,6	1.0	2.5	9.0
(VP held * talk) 3,6	1.1	2.4	9.5
(VP hold * conference) 3,6	3.5	5.1	17.0

# Cube Pruning

items are popped out-of-order

**solution:** keep a buffer of pop-ups

2.5 2.4 5.1

finally re-sort the buffer  
and return inorder:

2.4 2.5 5.1

(VP<sub>3,6</sub> held \* meeting)

(VP<sub>3,6</sub> held \* talk)

(VP<sub>3,6</sub> hold \* conference)

(PP<sub>1,3</sub> with \* Sharon)

(PP<sub>1,3</sub> along \* Sharon)

(PP<sub>1,3</sub> with \* Shalong)

	1.0	3.0	8.0
1.0	2.5	9.0	9.5
1.1	2.4	9.5	9.4
3.5	5.1	17.0	12.1