

Model Combination for Machine Translation

Google research

John DeNero, Shankar Kumar,
Ciprian Chelba, and Franz Josef Och

Motivation

Google

- ▶ A statistical machine translation model scores derivations
(log) model score sums **language model** and **translation model**

Motivation

- ▶ A statistical machine translation model scores derivations
(log) model score sums **language model** and **translation model**

$$\theta \cdot \phi(d) = \theta \cdot \left[\sum_{w \in \text{n-grams}(d)} \phi_{\text{LM}}(w) + \sum_{r \in \text{rules}(d)} \phi_{\text{TM}}(r) \right]$$

Motivation

- ▶ A statistical machine translation model scores derivations

(log) model score sums **language model** and **translation model**

$$\theta \cdot \phi(d) = \theta \cdot \left[\sum_{w \in \text{n-grams}(d)} \phi_{\text{LM}}(w) + \sum_{r \in \text{rules}(d)} \phi_{\text{TM}}(r) \right]$$

Motivation

- ▶ A statistical machine translation model scores derivations

(log) model score sums **language model** and **translation model**

$$\theta \cdot \phi(d) = \theta \cdot \left[\sum_{w \in \text{n-grams}(d)} \phi_{\text{LM}}(w) + \sum_{r \in \text{rules}(d)} \phi_{\text{TM}}(r) \right]$$

Motivation

- ▶ A statistical machine translation model scores derivations

(log) model score sums **language model** and **translation model**

$$\theta \cdot \phi(d) = \theta \cdot \left[\sum_{w \in \text{n-grams}(d)} \phi_{\text{LM}}(w) + \sum_{r \in \text{rules}(d)} \phi_{\text{TM}}(r) \right]$$

- ▶ We can improve over max-derivation decoding

Motivation

- ▶ A statistical machine translation model scores derivations

(log) model score sums **language model** and **translation model**

$$\theta \cdot \phi(d) = \theta \cdot \left[\sum_{w \in \text{n-grams}(d)} \phi_{\text{LM}}(w) + \sum_{r \in \text{rules}(d)} \phi_{\text{TM}}(r) \right]$$

- ▶ We can improve over max-derivation decoding
 - Consensus decoding (e.g., minimum Bayes risk)

Motivation

- ▶ A statistical machine translation model scores derivations

(log) model score sums **language model** and **translation model**

$$\theta \cdot \phi(d) = \theta \cdot \left[\sum_{w \in \text{n-grams}(d)} \phi_{\text{LM}}(w) + \sum_{r \in \text{rules}(d)} \phi_{\text{TM}}(r) \right]$$

- ▶ We can improve over max-derivation decoding

- Consensus decoding (e.g., minimum Bayes risk)



Motivation

- ▶ A statistical machine translation model scores derivations

(log) model score sums **language model** and **translation model**

$$\theta \cdot \phi(d) = \theta \cdot \left[\sum_{w \in \text{n-grams}(d)} \phi_{\text{LM}}(w) + \sum_{r \in \text{rules}(d)} \phi_{\text{TM}}(r) \right]$$

- ▶ We can improve over max-derivation decoding

- Consensus decoding (e.g., minimum Bayes risk)
- System combination (e.g., confusion networks)



Motivation

- ▶ A statistical machine translation model scores derivations

(log) model score sums **language model** and **translation model**

$$\theta \cdot \phi(d) = \theta \cdot \left[\sum_{w \in \text{n-grams}(d)} \phi_{\text{LM}}(w) + \sum_{r \in \text{rules}(d)} \phi_{\text{TM}}(r) \right]$$

- ▶ We can improve over max-derivation decoding

- Consensus decoding (e.g., minimum Bayes risk)
- System combination (e.g., confusion networks)



Motivation

- ▶ A statistical machine translation model scores derivations
(log) model score sums **language model** and **translation model**

$$\theta \cdot \phi(d) = \theta \cdot \left[\sum_{w \in \text{n-grams}(d)} \phi_{\text{LM}}(w) + \sum_{r \in \text{rules}(d)} \phi_{\text{TM}}(r) \right]$$

- ▶ We can improve over max-derivation decoding
 - Consensus decoding (e.g., minimum Bayes risk)
 - System combination (e.g., confusion networks)



- ▶ In this work, we develop a technique that integrates both

Consensus Decoding

- ▶ Derivation scores can be interpreted as probabilities

$$P(d|f) = \frac{\exp(\theta \cdot \phi(d))}{\sum_{d'} \exp(\theta \cdot \phi(d'))}$$

Consensus Decoding

Google

- ▶ Derivation scores can be interpreted as probabilities

$$P(d|f) = \frac{\exp(\theta \cdot \phi(d))}{\sum_{d'} \exp(\theta \cdot \phi(d'))}$$

- ▶ We can query this distribution for:

Consensus Decoding



- ▶ Derivation scores can be interpreted as probabilities

$$P(d|f) = \frac{\exp(\theta \cdot \phi(d))}{\sum_{d'} \exp(\theta \cdot \phi(d'))}$$

- ▶ We can query this distribution for:

Whole translations [Blunsom et al., '08]:

$$e_{\text{MAX-TRANS}} = \arg \max_e \sum_{d: \sigma_e(d)=e} P(d|f)$$

Consensus Decoding

- ▶ Derivation scores can be interpreted as probabilities

$$P(d|f) = \frac{\exp(\theta \cdot \phi(d))}{\sum_{d'} \exp(\theta \cdot \phi(d'))}$$

- ▶ We can query this distribution for:

Whole translations [Blunsom et al., '08]:

$$e_{\text{MAX-TRANS}} = \arg \max_e \sum_{d: \sigma_e(d)=e} P(d|f)$$

N-gram overlap [Kumar and Byrne, '04]:

$$e_{\text{MBR}} = \arg \max_e \mathbb{E} [\text{BLEU}(e; \sigma_e(d))]$$

Consensus Decoding

- ▶ Derivation scores can be interpreted as probabilities

$$P(d|f) = \frac{\exp(\theta \cdot \phi(d))}{\sum_{d'} \exp(\theta \cdot \phi(d'))}$$

- ▶ We can query this distribution for:

Whole translations [Blunsom et al., '08]:

$$e_{\text{MAX-TRANS}} = \arg \max_e \sum_{d: \sigma_e(d)=e} P(d|f)$$

N-gram overlap [Kumar and Byrne, '04]:

$$e_{\text{MBR}} = \arg \max_e \mathbb{E} [\text{BLEU}(e; \sigma_e(d))]$$



System Combination

Go.....oogle

- ▶ We often have multiple translation systems

System Combination

Google

- ▶ We often have multiple translation systems

el perro comí mi tarea

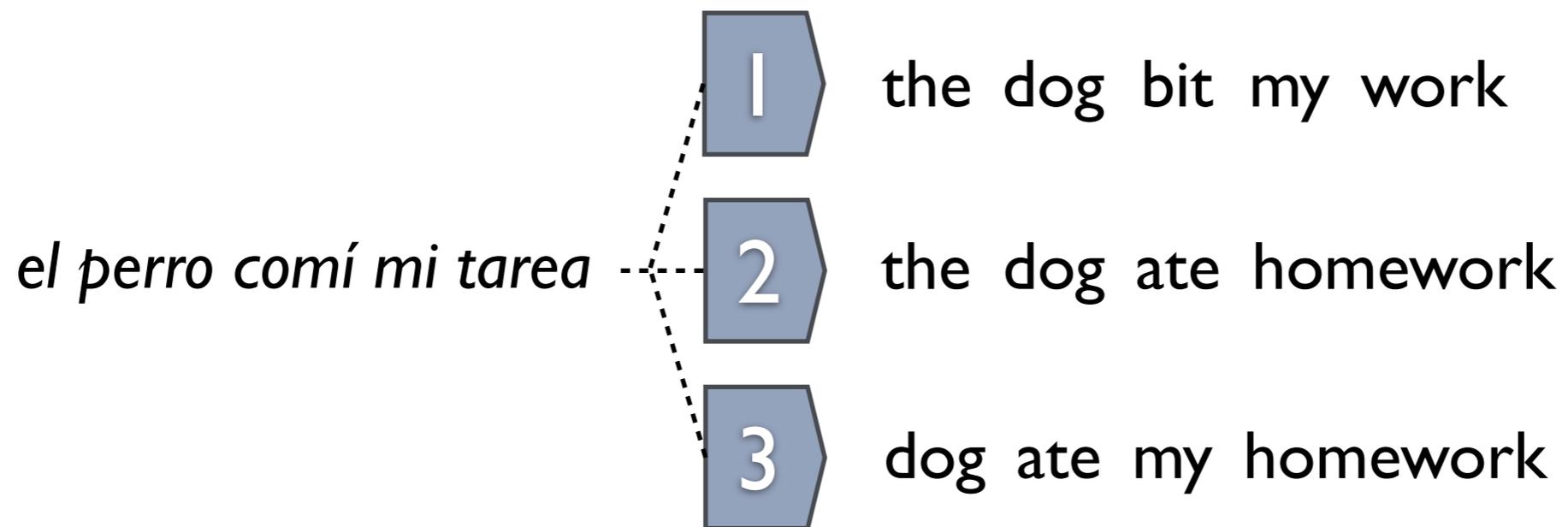
System Combination

- ▶ We often have multiple translation systems



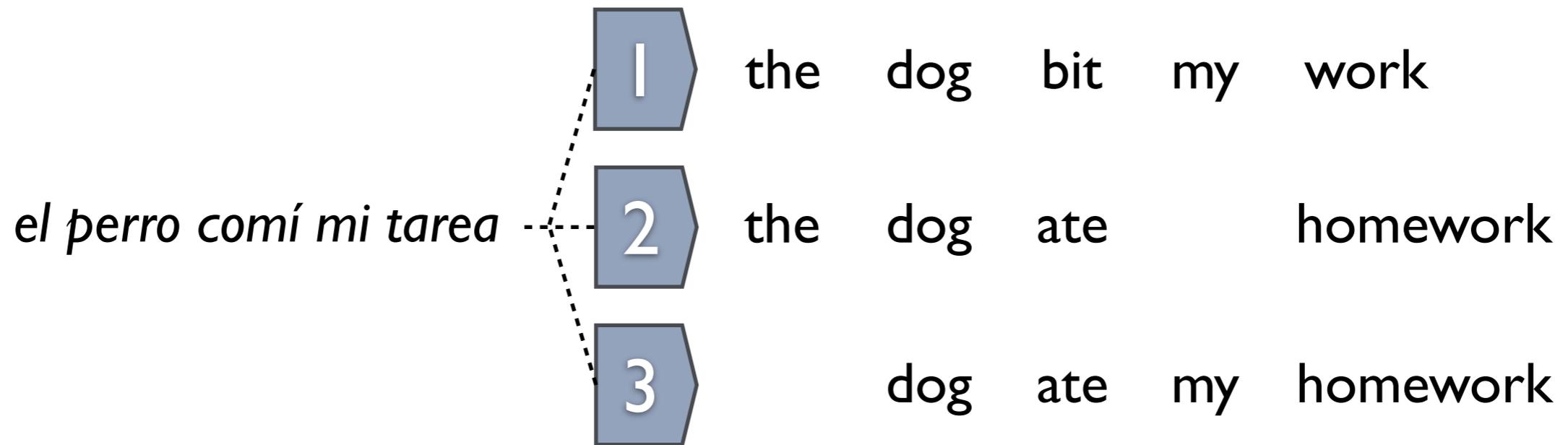
System Combination

- ▶ We often have multiple translation systems



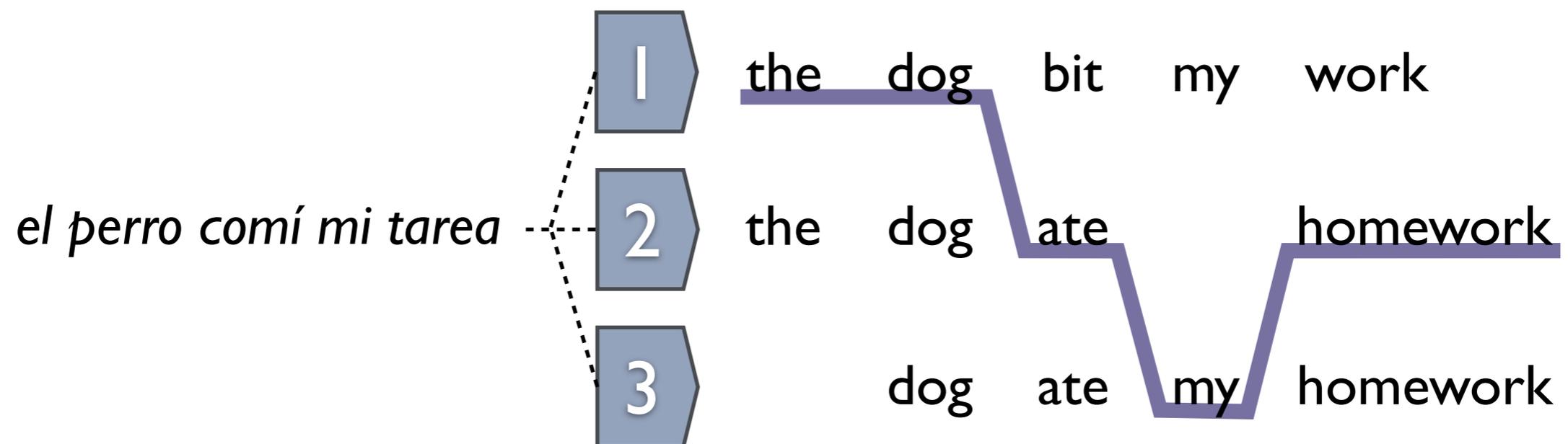
System Combination

- ▶ We often have multiple translation systems



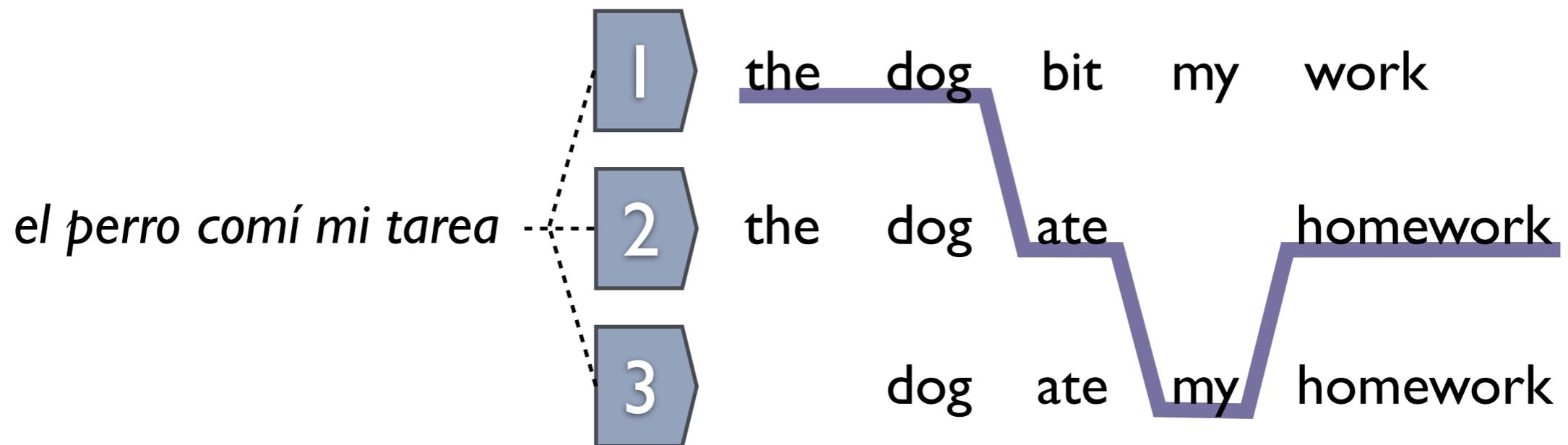
System Combination

- ▶ We often have multiple translation systems



System Combination

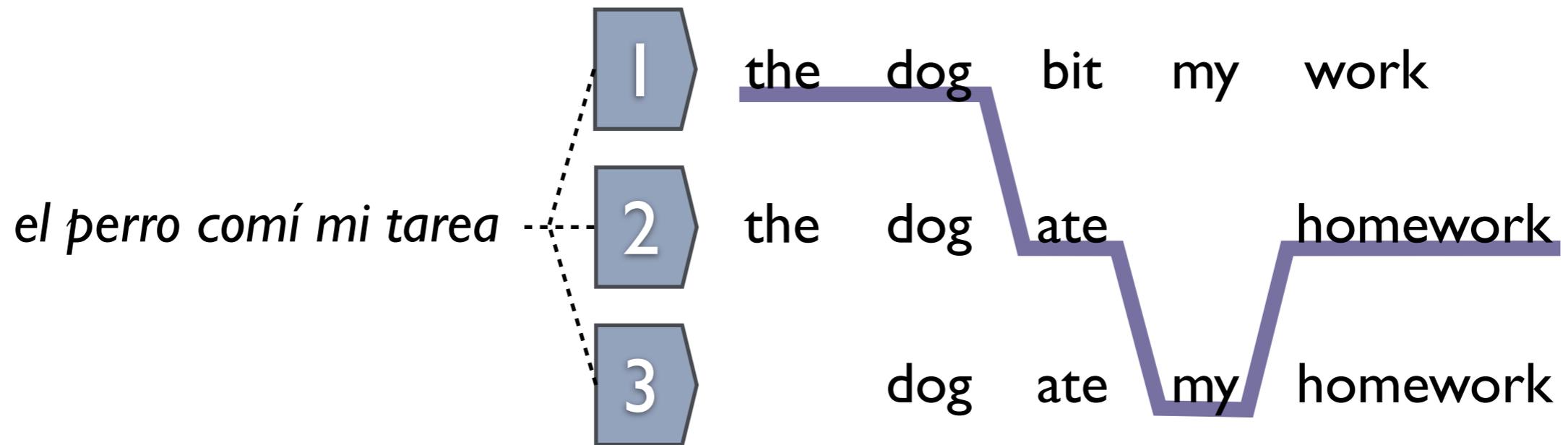
- ▶ We often have multiple translation systems



- ▶ Combiners assume little about systems

System Combination

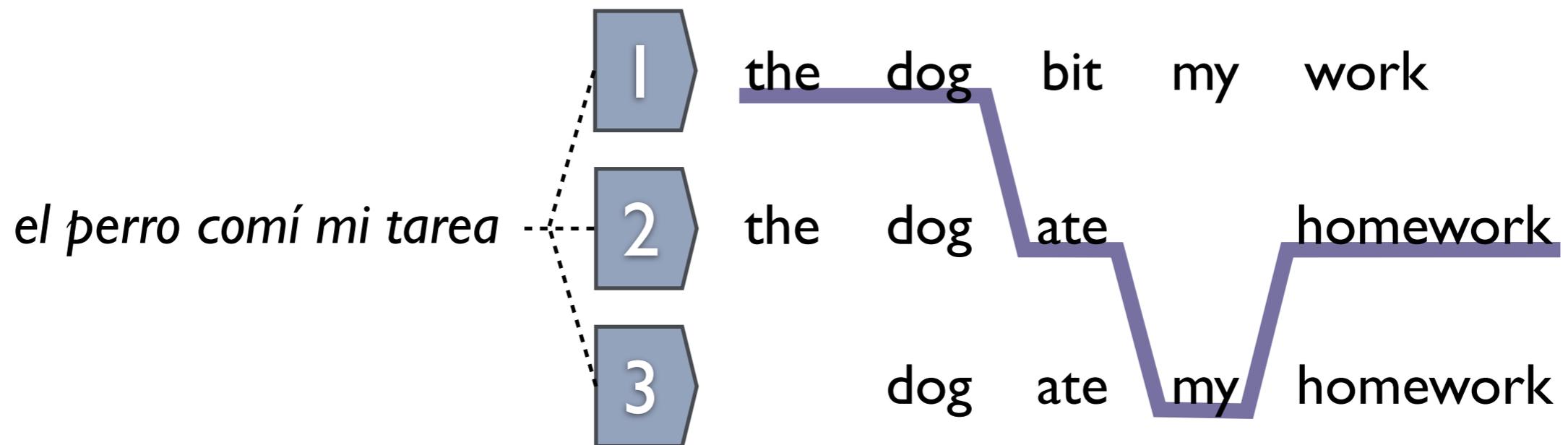
- ▶ We often have multiple translation systems



- ▶ Combiners assume little about systems
- ▶ Objectives similar to consensus decoding

System Combination

- ▶ We often have multiple translation systems



- ▶ Combiners assume little about systems
- ▶ Objectives similar to consensus decoding



Model Combination

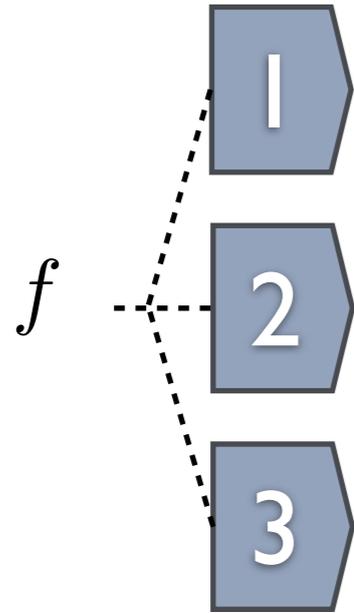


Model Combination

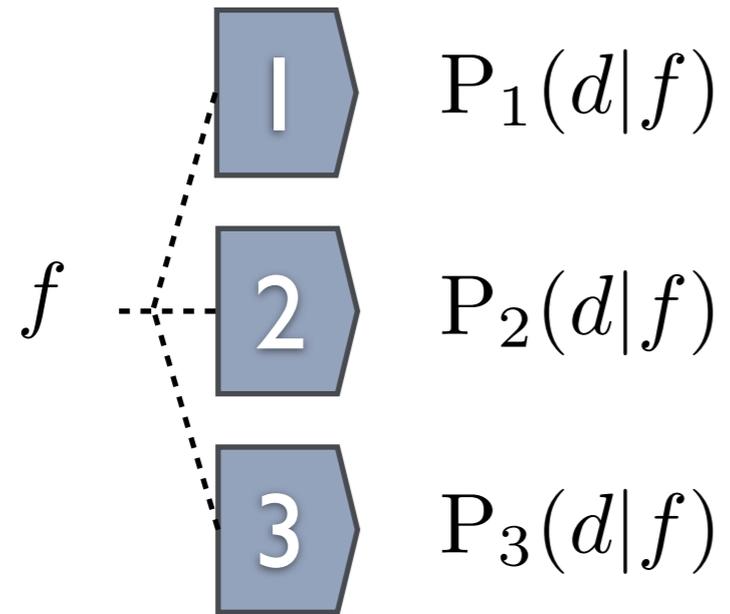
Google

f

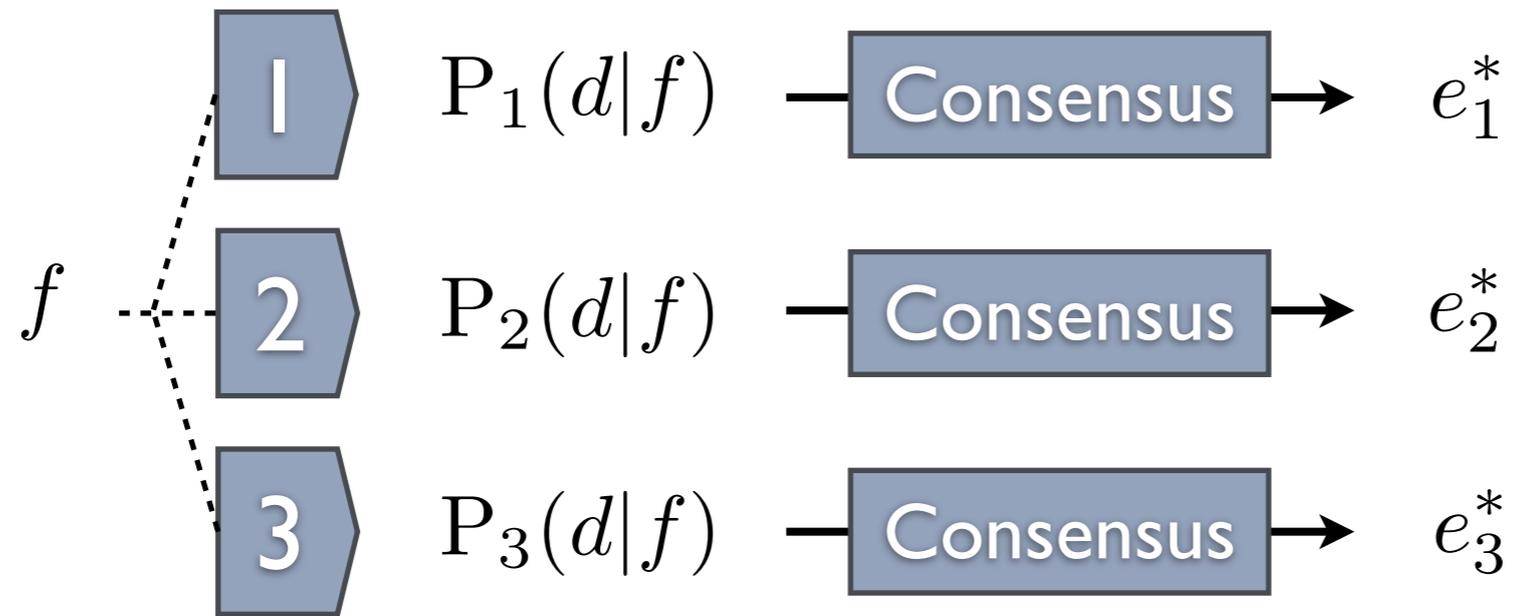
Model Combination



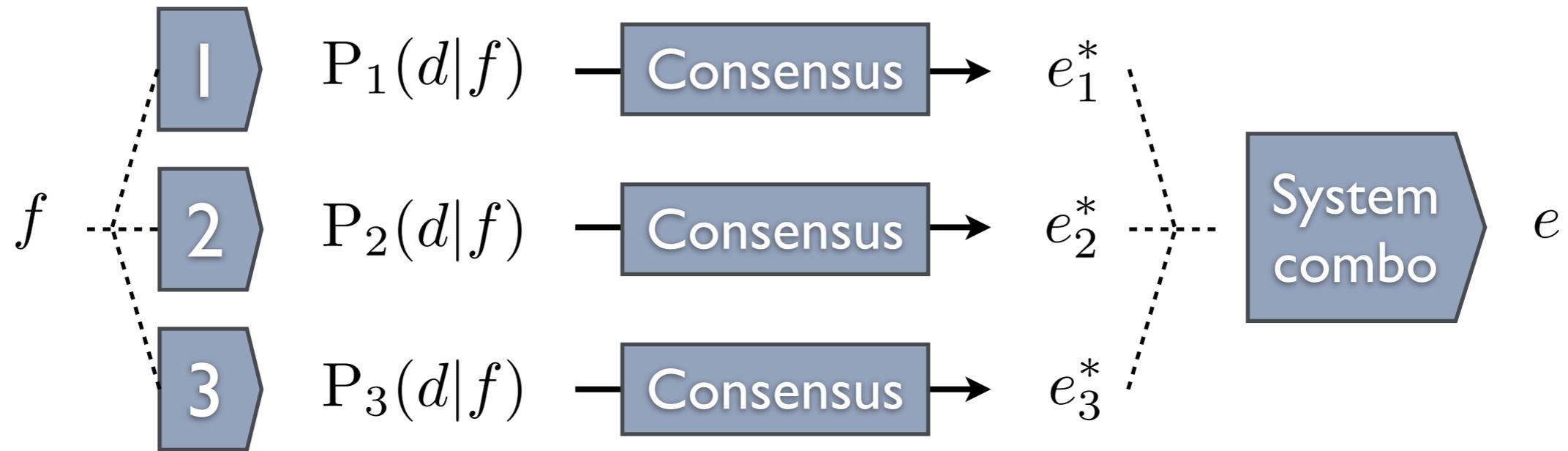
Model Combination



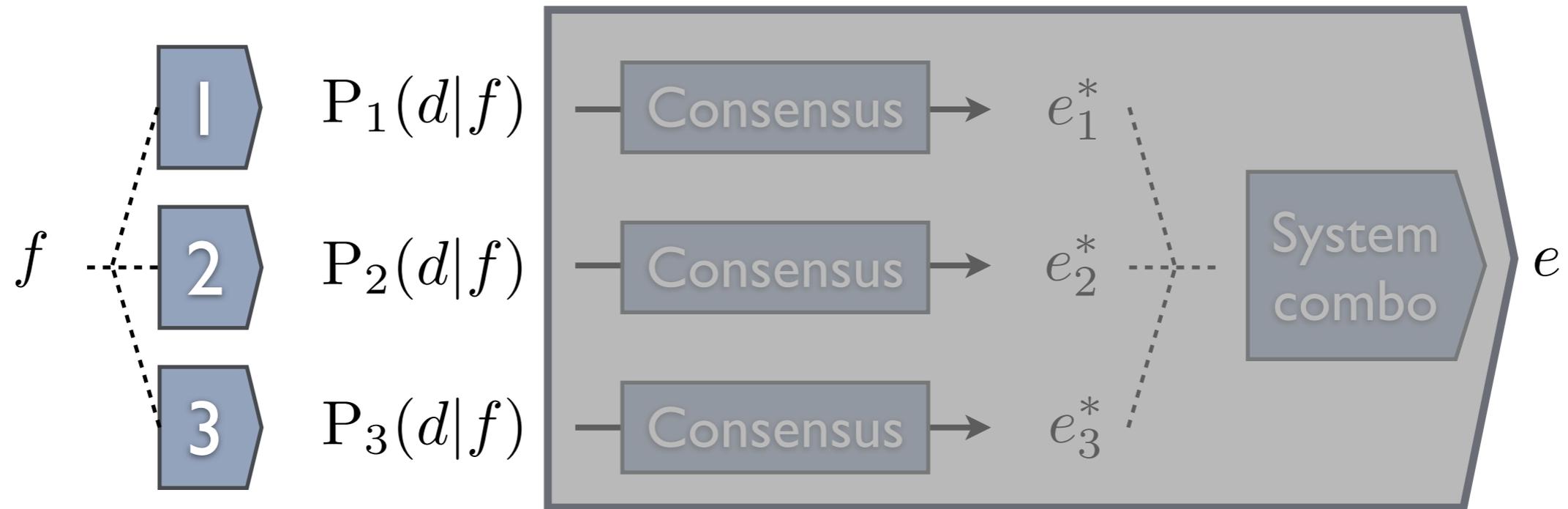
Model Combination



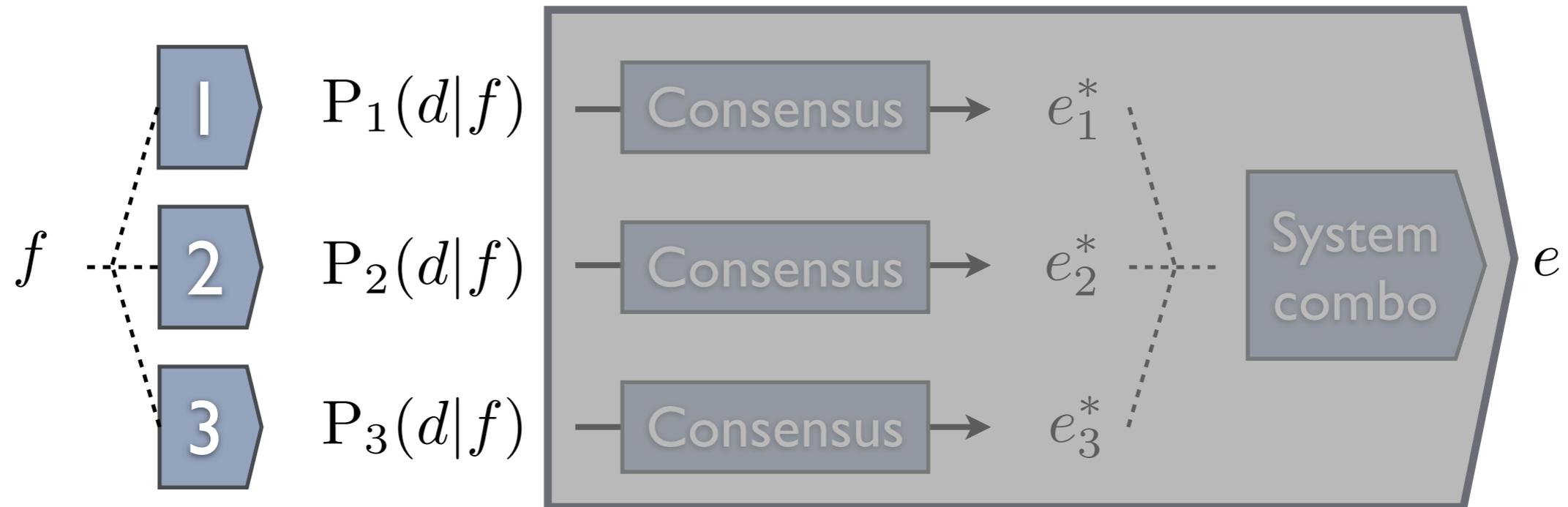
Model Combination



Model Combination

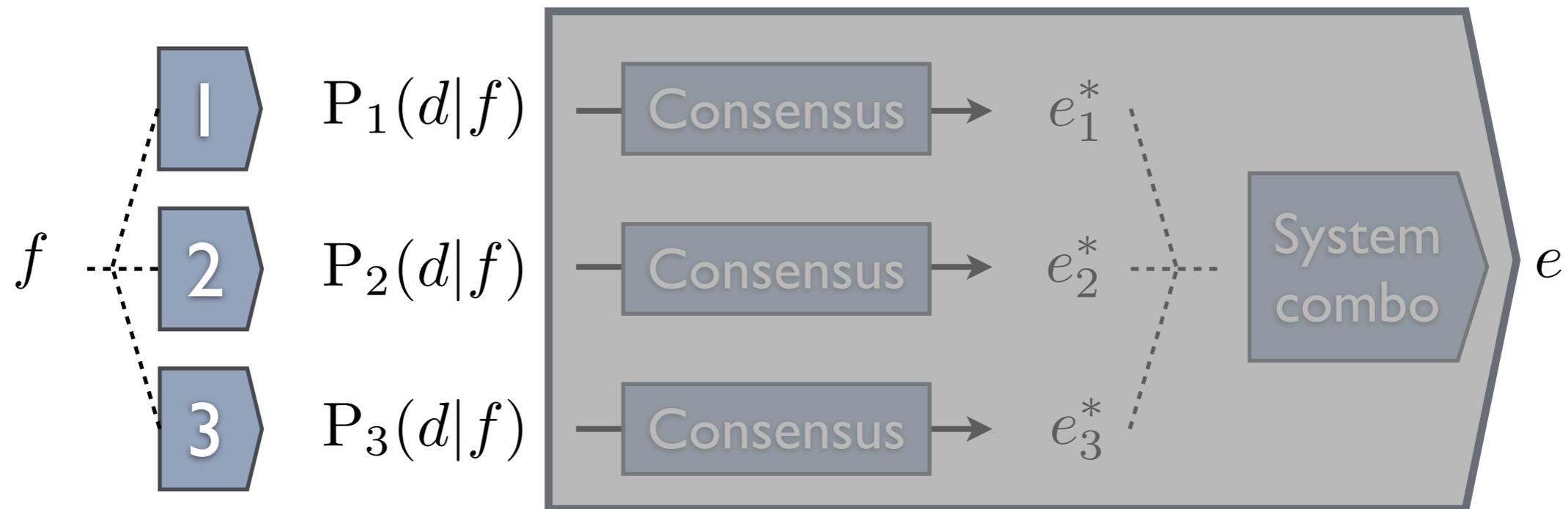


Model Combination



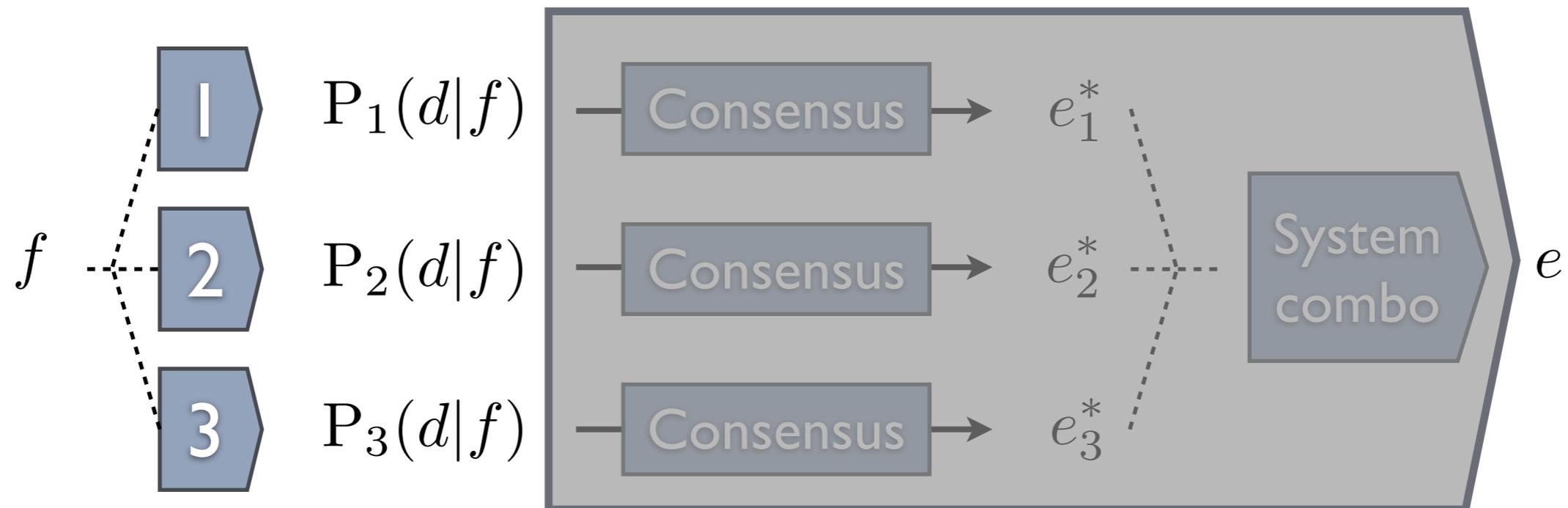
- ▶ Consensus decoding with multiple models

Model Combination



- ▶ Consensus decoding with multiple models
- ▶ Distribution-driven approach to system combination

Model Combination



- ▶ Consensus decoding with multiple models
- ▶ Distribution-driven approach to system combination
- ▶ Unifies consensus and combination objectives

Outline



Consensus decoding review

Our model combination technique

Comparison to system combination

Outline

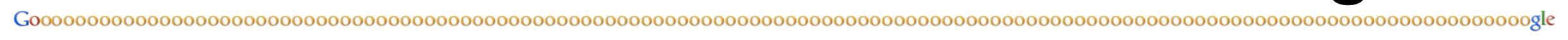
Consensus decoding review



Our model combination technique

Comparison to system combination

Forest-Based Consensus Decoding

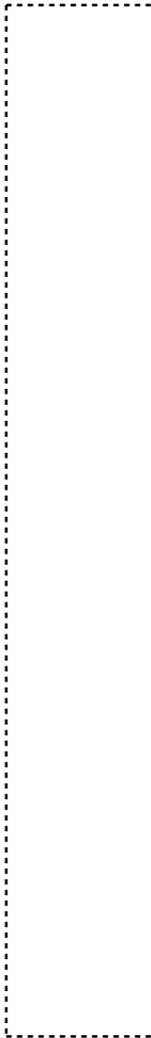


Forest-Based Consensus Decoding

Google

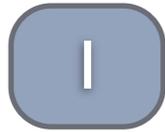


Build a forest that encodes the model posterior

$$P(d|f) =$$


Forest-Based Consensus Decoding

Google



Build a forest that encodes the model posterior

$$P(d|f) =$$

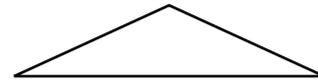


I saw



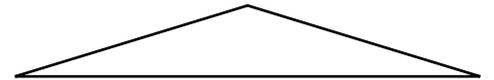
Yo vi

the man



al hombre

with {a,the} telescope



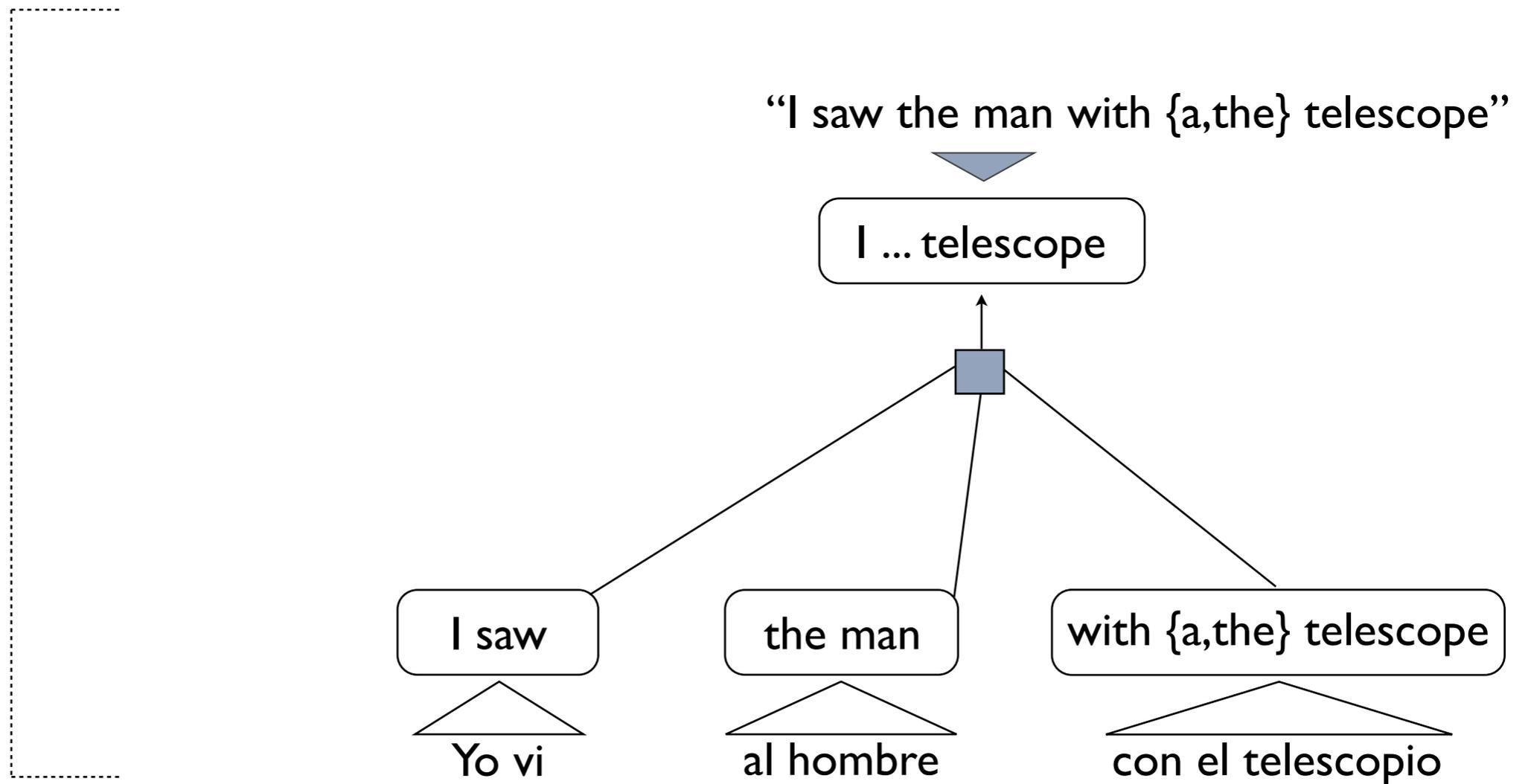
con el telescopio

Forest-Based Consensus Decoding

Google

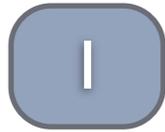
- I Build a forest that encodes the model posterior

$$P(d|f) =$$



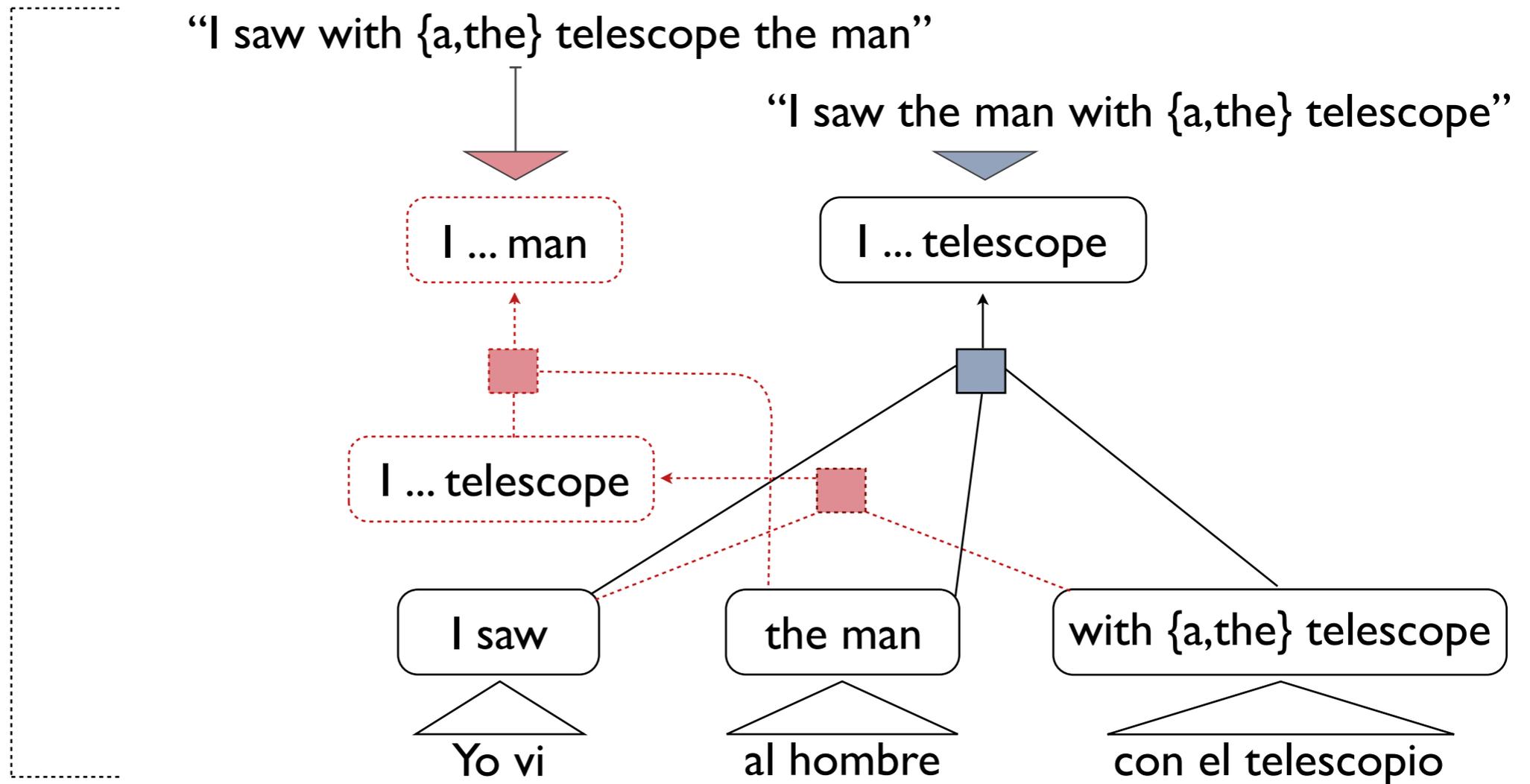
Forest-Based Consensus Decoding

Google



Build a forest that encodes the model posterior

$$P(d|f) =$$

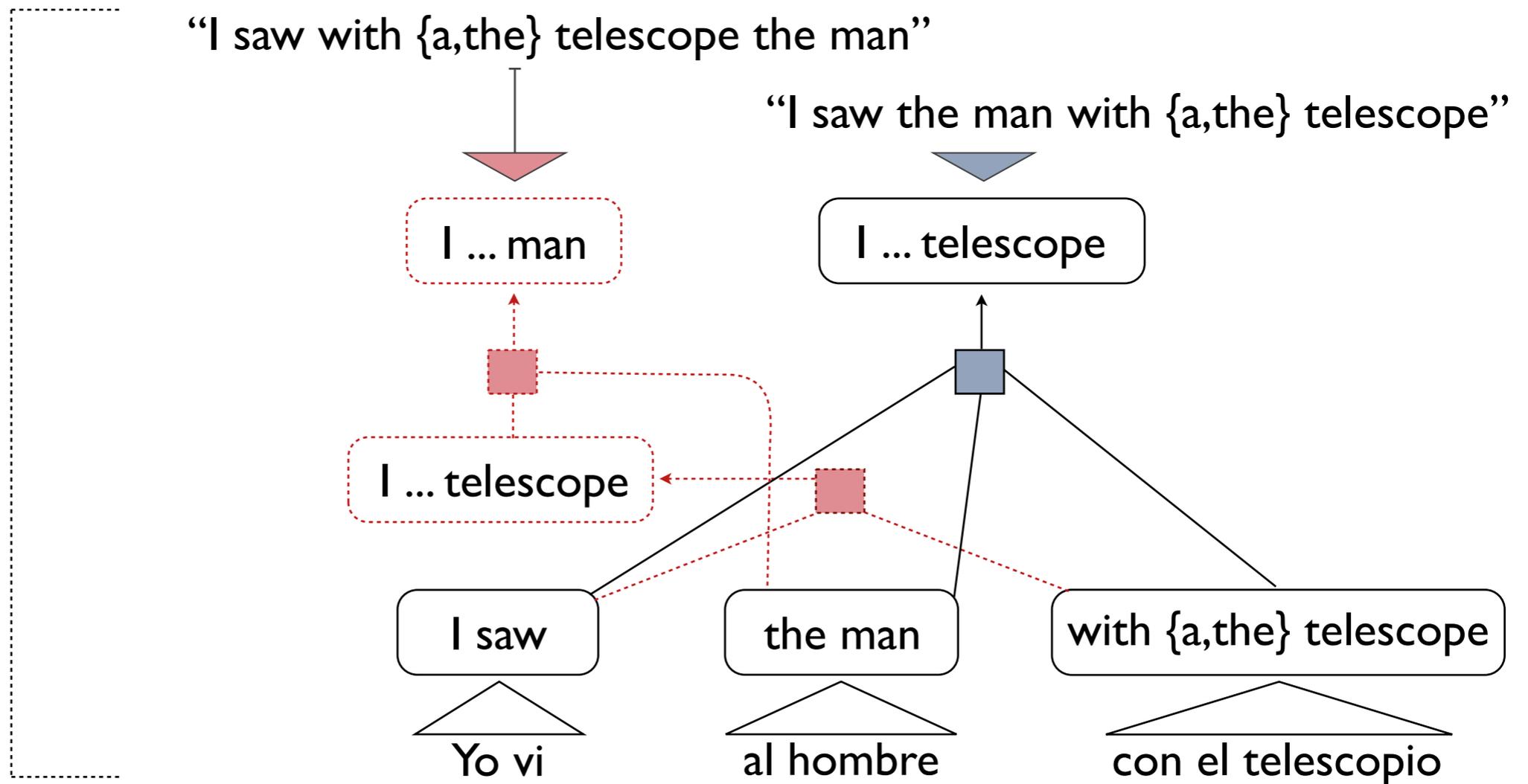


Forest-Based Consensus Decoding

Google

- 1 Build a forest that encodes the model posterior
- 2 Compute *n-gram statistics* from the posterior

$$P(d|f) =$$

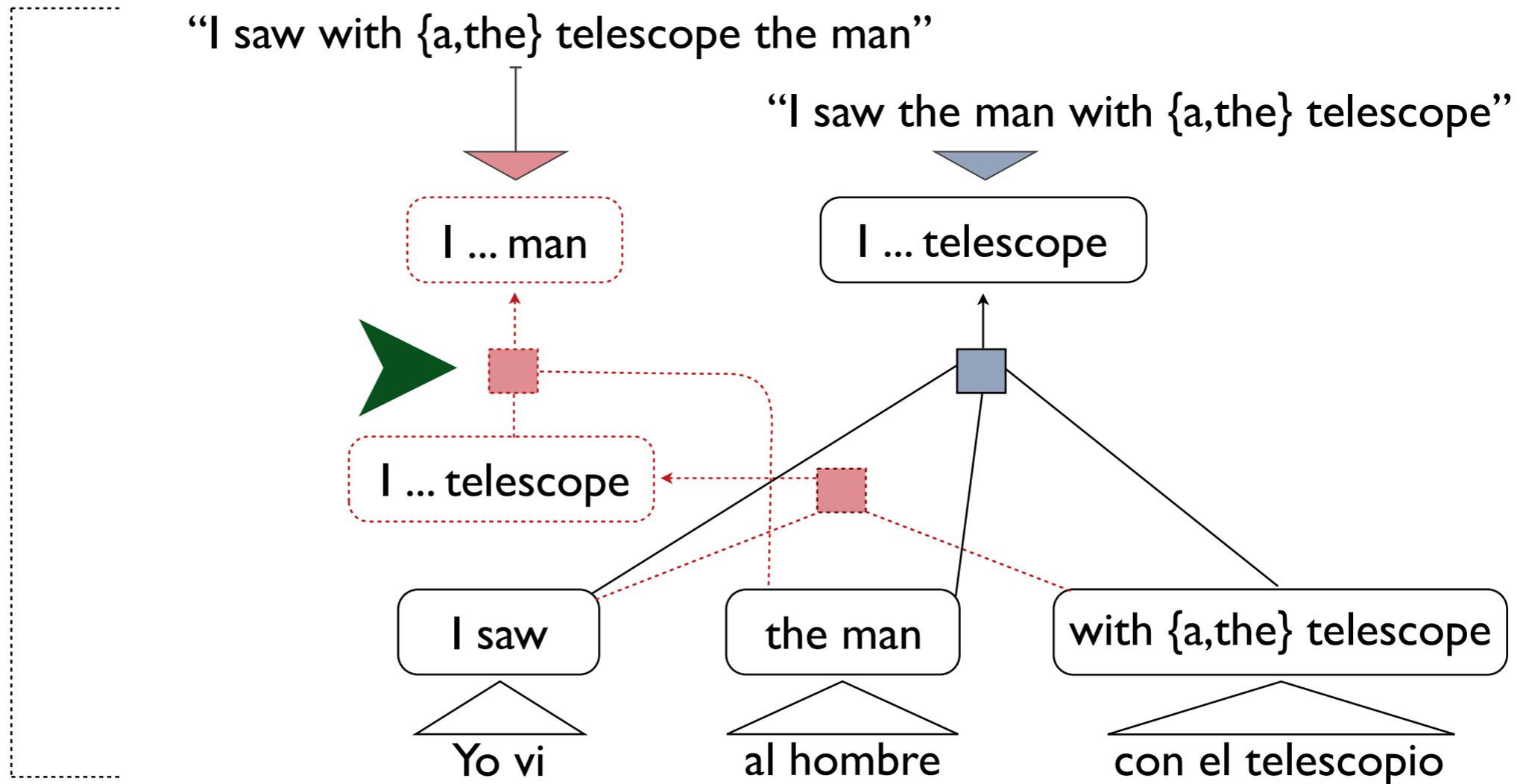


Forest-Based Consensus Decoding

Google

- 1 Build a forest that encodes the model posterior
- 2 Compute *n-gram statistics* from the posterior

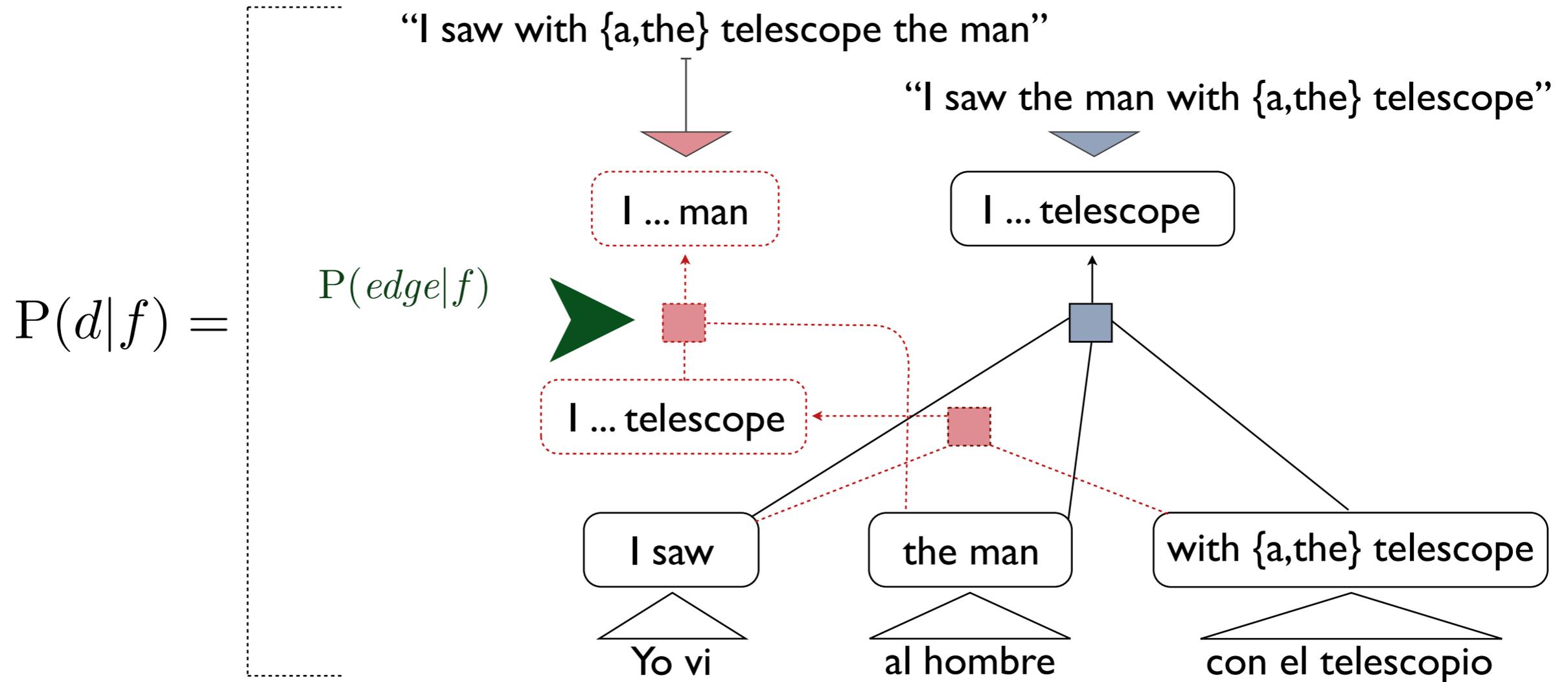
$$P(d|f) =$$



Forest-Based Consensus Decoding

Google

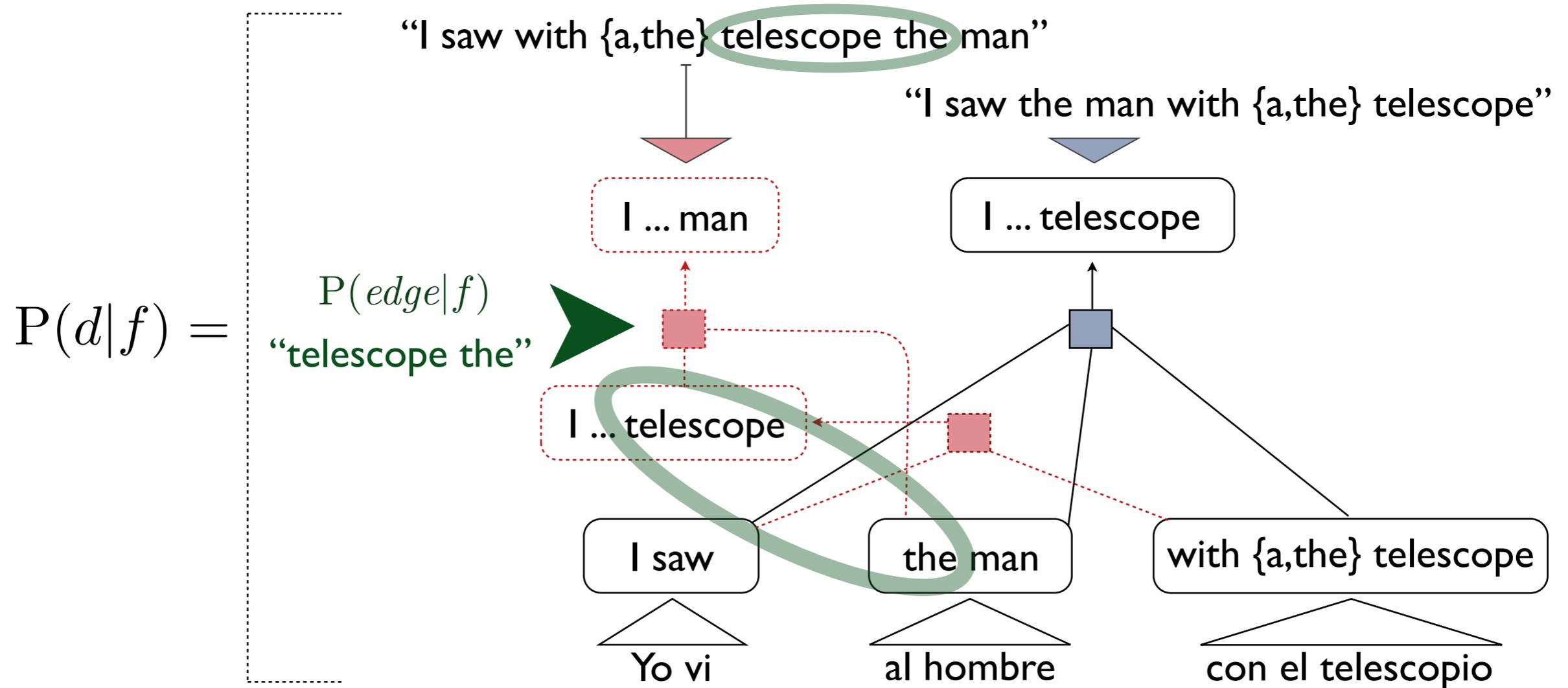
- 1 Build a forest that encodes the model posterior
- 2 Compute *n-gram statistics* from the posterior



Forest-Based Consensus Decoding

Google

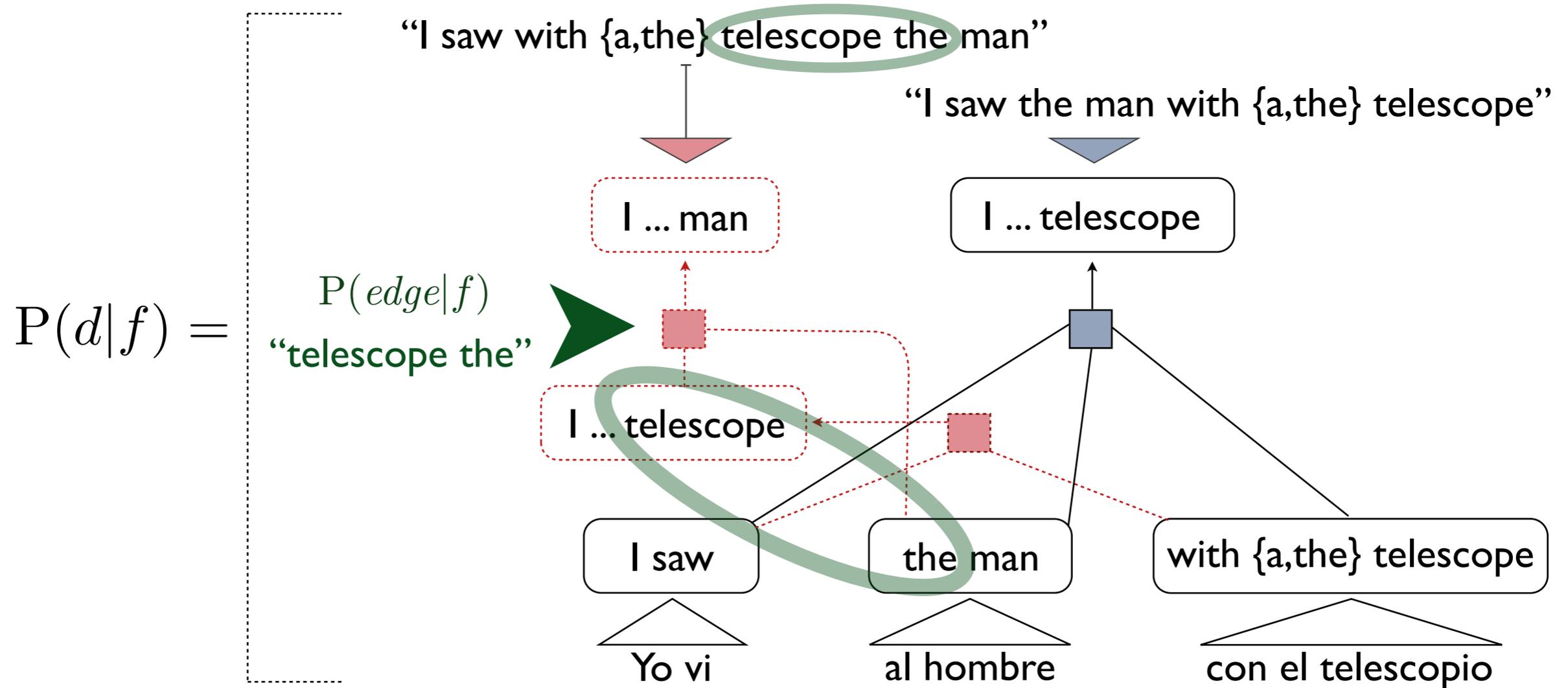
- 1 Build a forest that encodes the model posterior
- 2 Compute *n-gram statistics* from the posterior



Forest-Based Consensus Decoding

Google

- 1 Build a forest that encodes the model posterior
- 2 Compute *n-gram statistics* from the posterior
- 3 Optimize a *consensus objective* using these statistics



Types of Efficient Consensus Techniques



<p><i>Lattice Minimum Bayes-Risk Decoding</i></p> <p>[Tromble et al., '08]</p>	

Types of Efficient Consensus Techniques

Google

<i>Lattice Minimum Bayes-Risk Decoding</i> [Tromble et al., '08]	

Posteriors

Expected counts

N-gram Statistics

Types of Efficient Consensus Techniques

Google

Consensus Objective	Learned		
	Fixed	<i>Lattice Minimum Bayes-Risk Decoding</i> [Tromble et al., '08]	
		Posteriors	Expected counts

N-gram Statistics

Types of Efficient Consensus Techniques

Google

Consensus Objective	Learned	<p><i>Minimum Bayes-Risk Decoding for Hypergraphs</i></p> <p>[Kumar et al., '09]</p>	<p><i>Variational Decoding for Machine Translation</i></p> <p>[Li et al., '09]</p>
	Fixed	<p><i>Lattice Minimum Bayes-Risk Decoding</i></p> <p>[Tromble et al., '08]</p>	<p><i>Fast Consensus Decoding</i></p> <p>[DeNero et al., '09]</p>
		Posteriors	Expected counts
<hr/> <p>N-gram Statistics</p>			

Review

Learned Objectives are Better than Fixed

Google

Learned Objectives are Better than Fixed

Google

$$C(d) = \sum_{n=1}^4 w^{(n)} \sum_{g \in n\text{-grams}} c(g, d) \cdot P(g|f) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} \theta \cdot \phi(d)$$

Learned Objectives are Better than Fixed

Google

N-gram statistics

$$C(d) = \sum_{n=1}^4 w^{(n)} \sum_{g \in n\text{-grams}} c(g, d) \cdot P(g|f) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} \theta \cdot \phi(d)$$

Learned Objectives are Better than Fixed

Google

N-gram statistics

Length

$$C(d) = \sum_{n=1}^4 w^{(n)} \sum_{g \in n\text{-grams}} c(g, d) \cdot P(g|f) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} \theta \cdot \phi(d)$$

Learned Objectives are Better than Fixed

Google

N-gram statistics

Length

Base model

$$C(d) = \sum_{n=1}^4 w^{(n)} \sum_{g \in n\text{-grams}} c(g, d) \cdot P(g|f) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} \theta \cdot \phi(d)$$

Learned Objectives are Better than Fixed

Google

$$C(d) = \sum_{n=1}^4 w^{(n)} \sum_{g \in n\text{-grams}} c(g, d) \cdot P(g|f) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} \theta \cdot \phi(d)$$

Objective parameters

Fixed: Choose w such that $C_w(d) \approx \mathbb{E} [\text{BLEU}(d)]$

Learned: Choose w to **maximize** $\text{BLEU} \left(\left\{ \arg \max_d C_w(d) \right\} ; \mathbf{e} \right)$
[Kumar et al., '09]

Learned Objectives are Better than Fixed

Google

$$C(d) = \sum_{n=1}^4 w^{(n)} \sum_{g \in n\text{-grams}} c(g, d) \cdot P(g|f) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} \theta \cdot \phi(d)$$

Objective parameters

Fixed: Choose w such that $C_w(d) \approx \mathbb{E} [\text{BLEU}(d)]$

References

Learned: Choose w to **maximize** $\text{BLEU} \left(\left\{ \arg \max_d C_w(d) \right\}; \mathbf{e} \right)$
 [Kumar et al., '09]

Learned Objectives are Better than Fixed

Google

$$C(d) = \sum_{n=1}^4 w^{(n)} \sum_{g \in n\text{-grams}} c(g, d) \cdot P(g|f) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} \theta \cdot \phi(d)$$

Objective parameters

Fixed: Choose w such that $C_w(d) \approx \mathbb{E} [\text{BLEU}(d)]$ *References*

Learned: Choose w to **maximize** $\text{BLEU} \left(\left\{ \arg \max_d C_w(d) \right\}; \mathbf{e} \right)$
[Kumar et al., '09]

Consensus performance versus max-derivation decoding (39 pairs)

- Increased test set BLEU by ≥ 0.2
- Decreased test set BLEU by ≥ 0.2

Learned Objectives are Better than Fixed

Google

$$C(d) = \sum_{n=1}^4 w^{(n)} \sum_{g \in n\text{-grams}} c(g, d) \cdot P(g|f) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} \theta \cdot \phi(d)$$

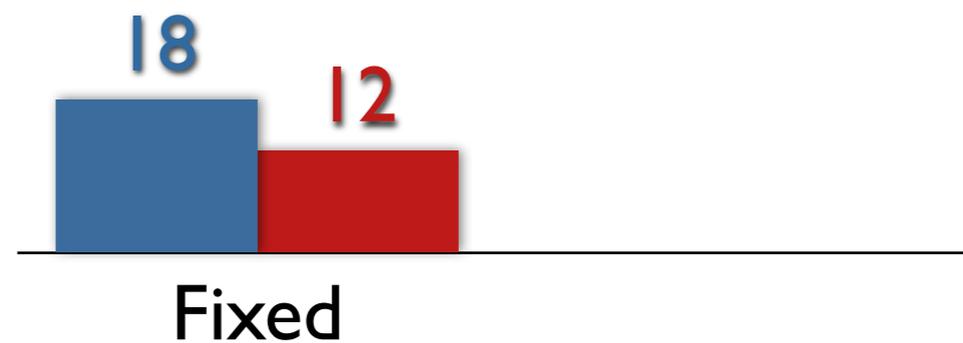
Objective parameters

Fixed: Choose w such that $C_w(d) \approx \mathbb{E} [\text{BLEU}(d)]$ References

Learned: Choose w to **maximize** $\text{BLEU} \left(\left\{ \arg \max_d C_w(d) \right\}; e \right)$
[Kumar et al., '09]

Consensus performance versus max-derivation decoding (39 pairs)

- Increased test set BLEU by ≥ 0.2
- Decreased test set BLEU by ≥ 0.2



Learned Objectives are Better than Fixed

Google

$$C(d) = \sum_{n=1}^4 w^{(n)} \sum_{g \in n\text{-grams}} c(g, d) \cdot P(g|f) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} \theta \cdot \phi(d)$$

Objective parameters

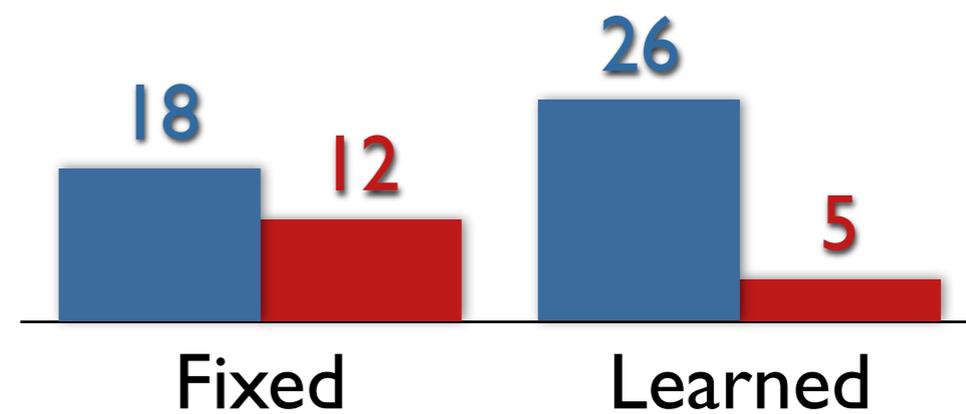
Fixed: Choose w such that $C_w(d) \approx \mathbb{E} [\text{BLEU}(d)]$

References

Learned: Choose w to **maximize** $\text{BLEU} \left(\left\{ \arg \max_d C_w(d) \right\}; e \right)$
[Kumar et al., '09]

Consensus performance versus max-derivation decoding (39 pairs)

- Increased test set BLEU by ≥ 0.2
- Decreased test set BLEU by ≥ 0.2



N-gram Posteriors can also be Computed Quickly

Google

N-gram Posteriors can also be Computed Quickly

Google

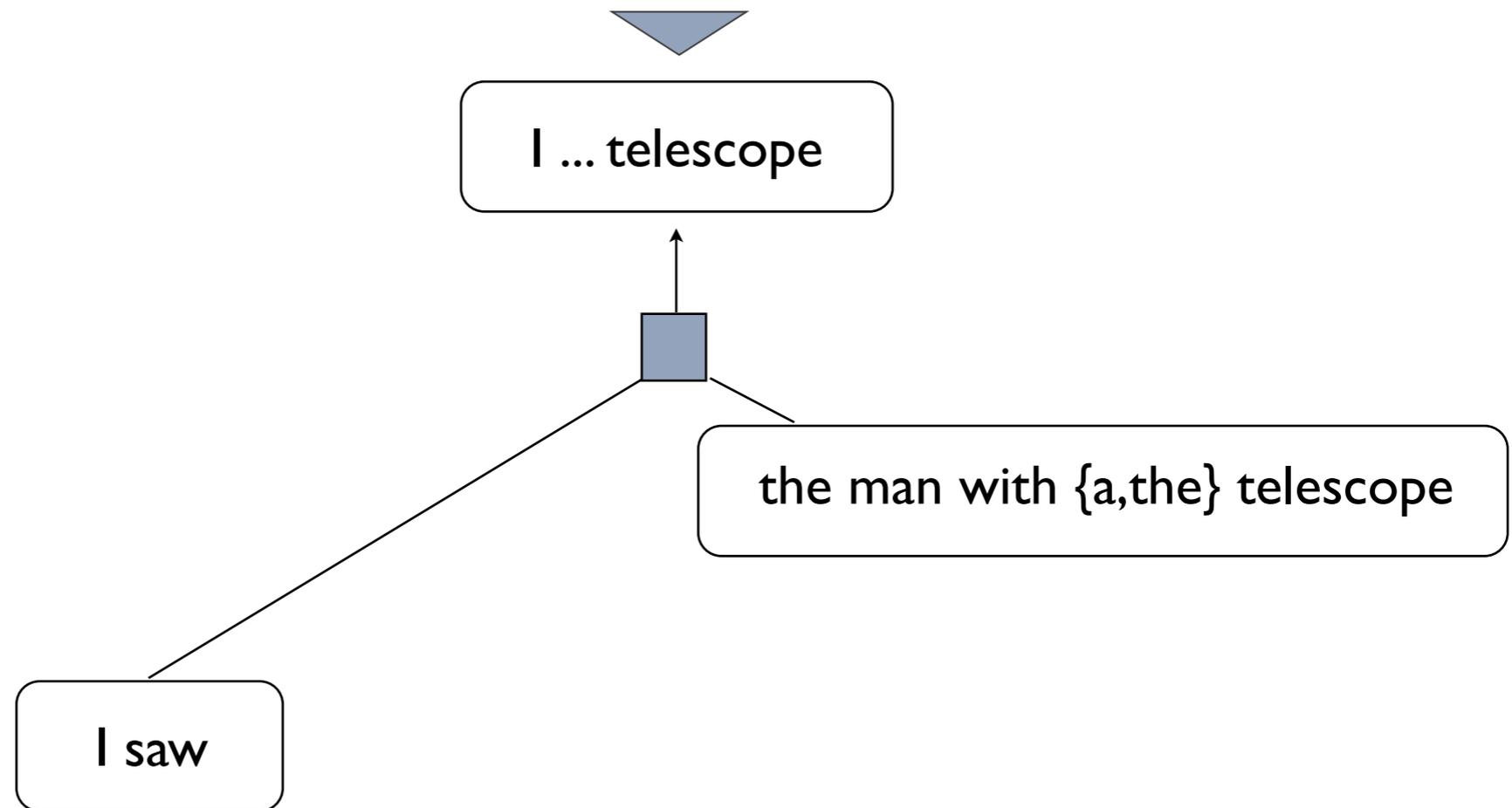
$$P(g|f) = 1.0 - P(\bar{g}|f)$$

N-gram Posteriors can also be Computed Quickly

Google

$$P(g|f) = 1.0 - P(\bar{g}|f)$$

“I saw the man with {a,the} telescope”

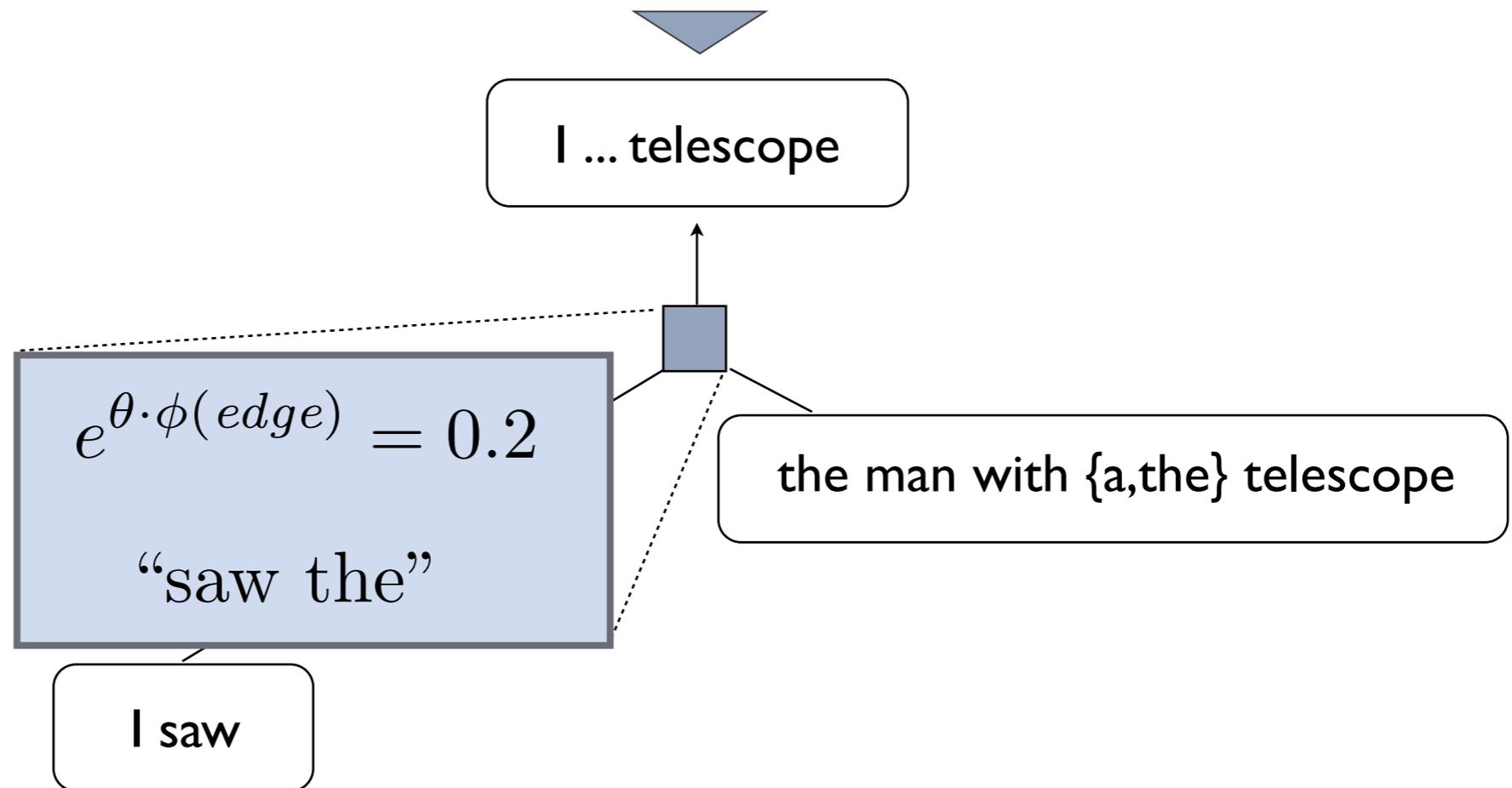


N-gram Posteriors can also be Computed Quickly

Google

$$P(g|f) = 1.0 - P(\bar{g}|f)$$

“I saw the man with {a,the} telescope”



N-gram Posteriors can also be Computed Quickly

Google

$$P(g|f) = 1.0 - P(\bar{g}|f)$$

“I saw the man with {a,the} telescope”

I ... telescope

$$e^{\theta \cdot \phi(\text{edge})} = 0.2$$

“saw the”

the man with {a,the} telescope

Inside: 0.1

“I saw” : 0

“with the” : 0.1

...

I saw

N-gram Posteriors can also be Computed Quickly

Google

$$P(g|f) = 1.0 - P(\bar{g}|f)$$

“I saw the man with {a,the} telescope”

I ... telescope

$$e^{\theta \cdot \phi(\text{edge})} = 0.2$$

“saw the”

the man with {a,the} telescope

Inside: 0.1

“I saw” : 0

“with the” : 0.1

...

I saw

Derivations that don't contain “with the”

N-gram Posteriors can also be Computed Quickly

Google

$$P(g|f) = 1.0 - P(\bar{g}|f)$$

“I saw the man with {a,the} telescope”

I ... telescope

$$e^{\theta \cdot \phi(\text{edge})} = 0.2$$

“saw the”

the man with {a,the} telescope

Inside: 0.1

“I saw” : 0

“with the” : 0.1

...

I saw

Derivations that don't contain “with the”

Inside: 0.5

“with a” : 0.3

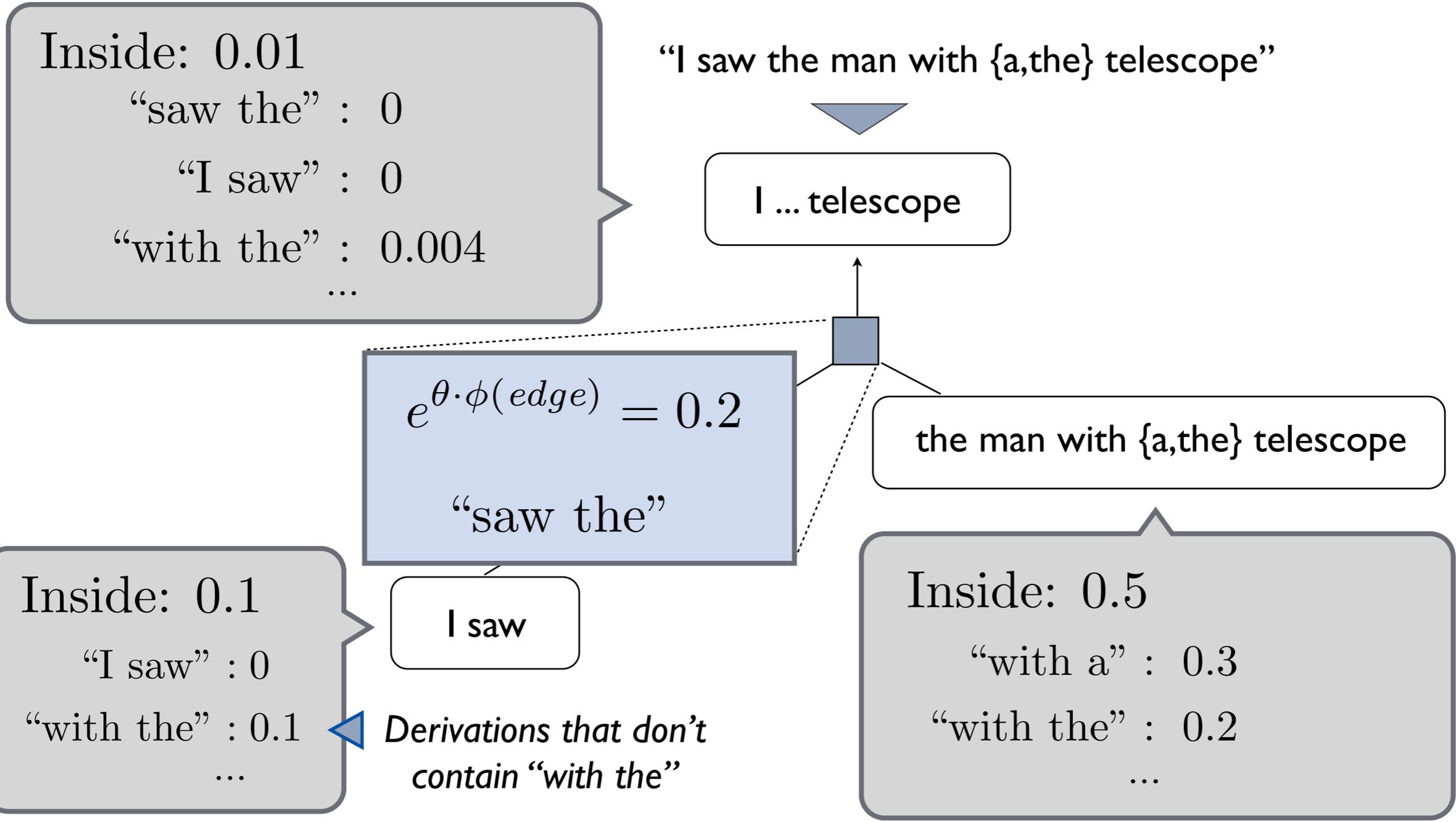
“with the” : 0.2

...

N-gram Posteriors can also be Computed Quickly



$$P(g|f) = 1.0 - P(\bar{g}|f)$$



N-gram Posteriors can also be Computed Quickly



$$P(g|f) = 1.0 - P(\bar{g}|f)$$

“I saw the man with {a,the} telescope”

I ... telescope

the man with {a,the} telescope

Inside: 0.5

“with a” : 0.3

“with the” : 0.2

...

“saw the”

I saw

Derivations that don't contain “with the”

$$e^{\theta \cdot \phi(\text{edge})} = 0.2$$

Inside: 0.01

“saw the” : 0

“I saw” : 0

“with the” : 0.004

...

Inside: 0.1

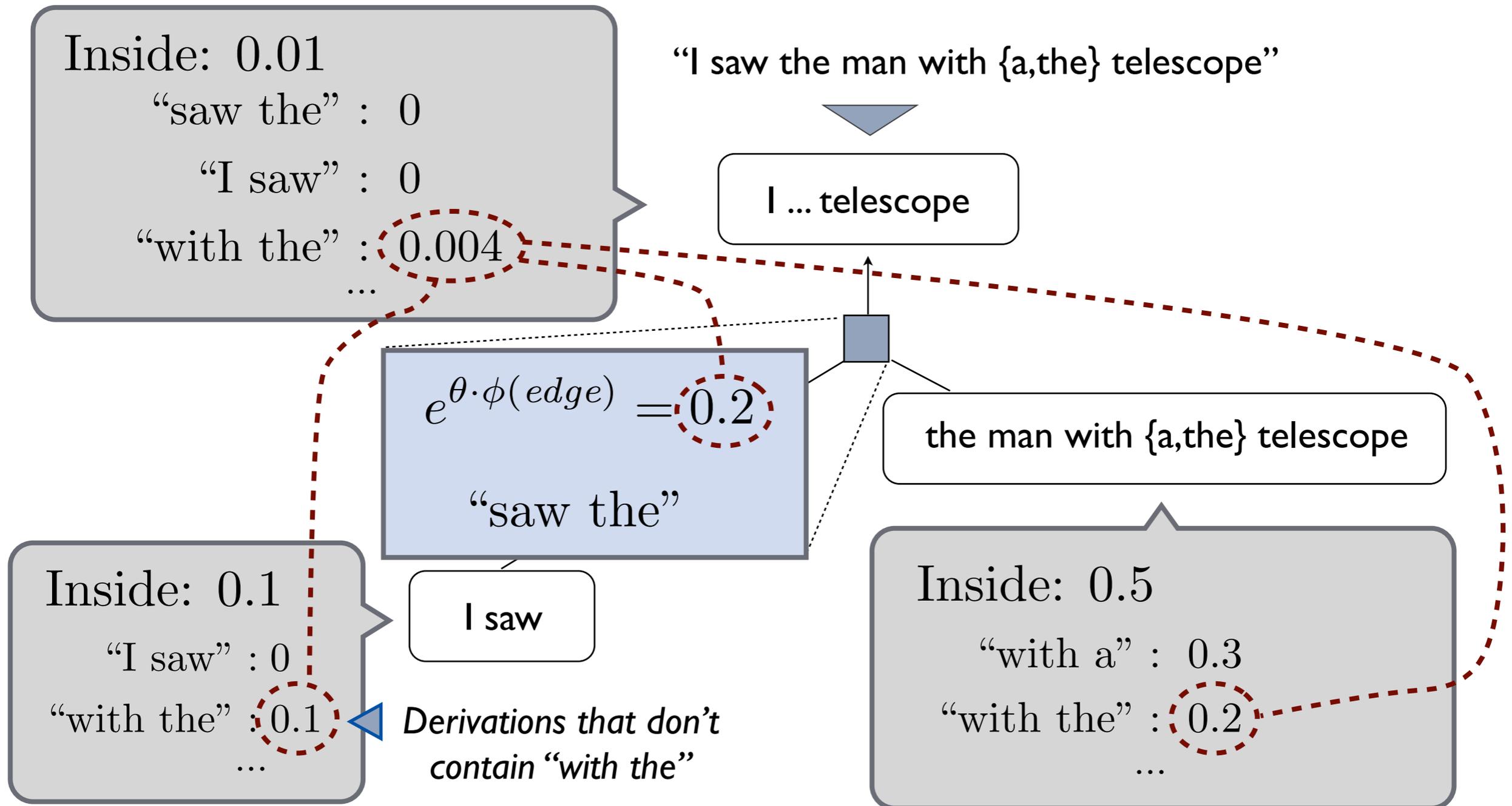
“I saw” : 0

“with the” : 0.1

...

N-gram Posteriors can also be Computed Quickly

$$P(g|f) = 1.0 - P(\bar{g}|f)$$

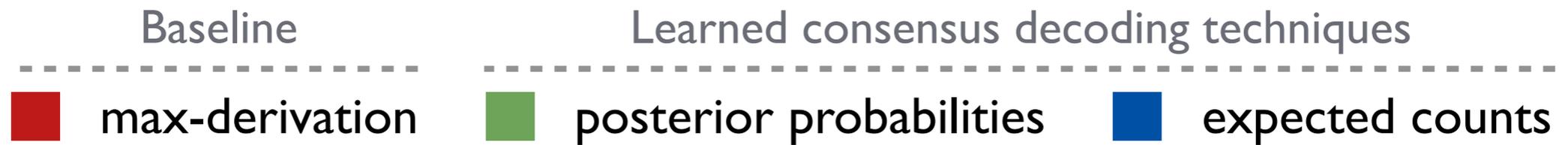


Audience challenge: What semiring computes n-gram posteriors?

Results for Learned Consensus Decoding

Google

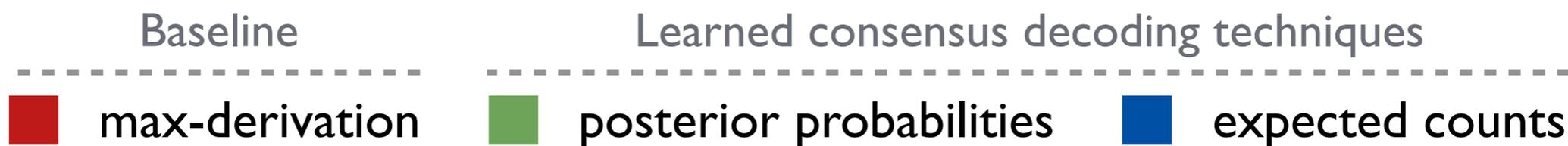
Constrained data track of the 2008 NIST MT task



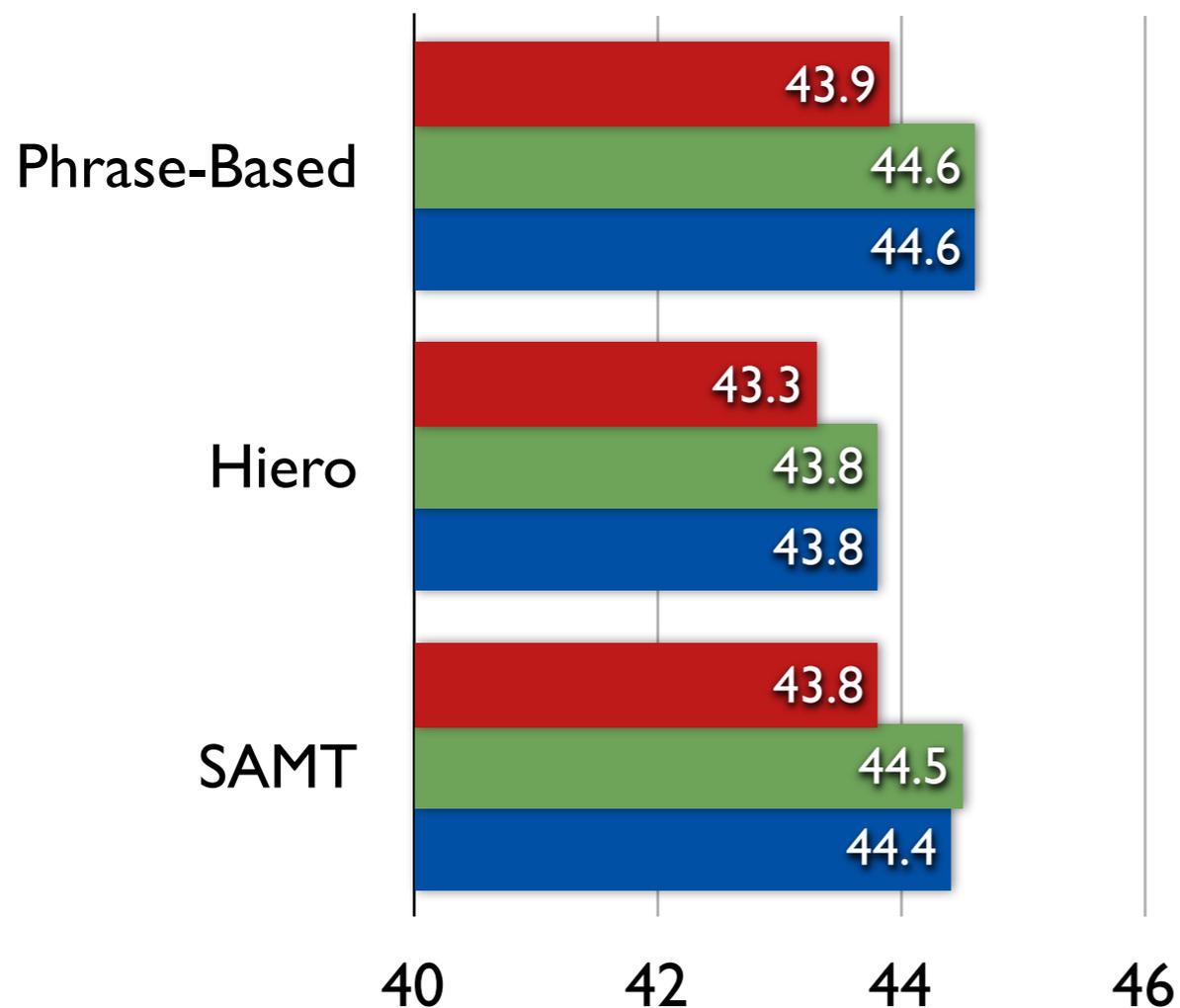
Results for Learned Consensus Decoding

Google

Constrained data track of the 2008 NIST MT task



Arabic-to-English



Results for Learned Consensus Decoding

Google

Constrained data track of the 2008 NIST MT task

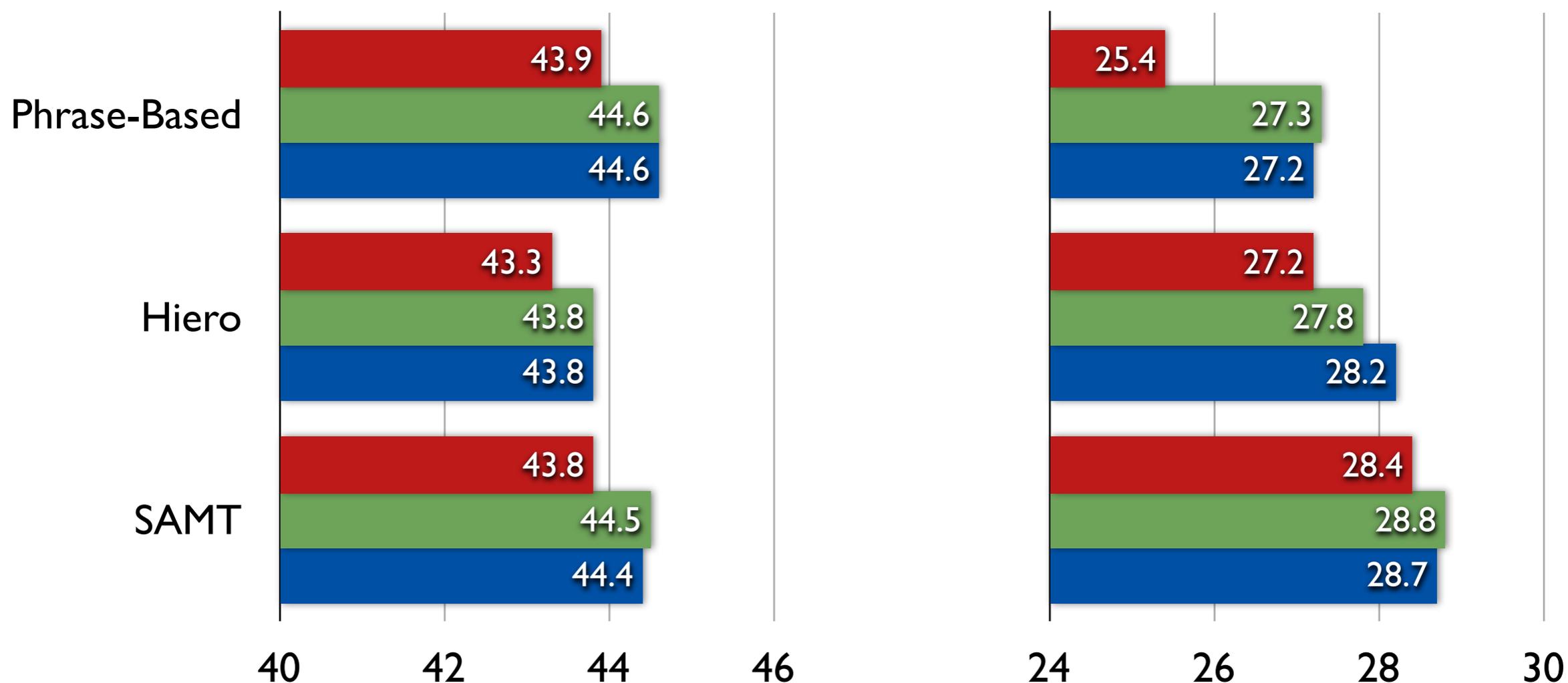
Baseline

Learned consensus decoding techniques

max-derivation posterior probabilities expected counts

Arabic-to-English

Chinese-to-English



Outline



Consensus decoding review

Our model combination technique

Comparison to system combination

Outline

Consensus decoding review

Our model combination technique



**All in
One Slide!**

Comparison to system combination

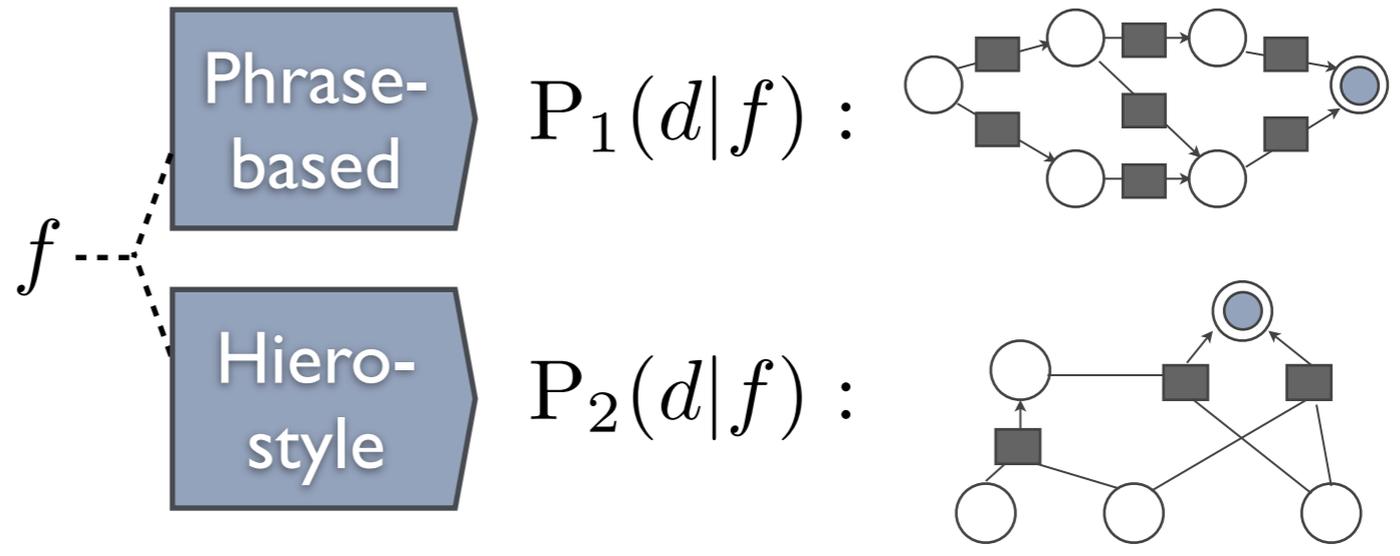
Extending Consensus Decoding to Multiple Models

Google

I. Build posterior forests

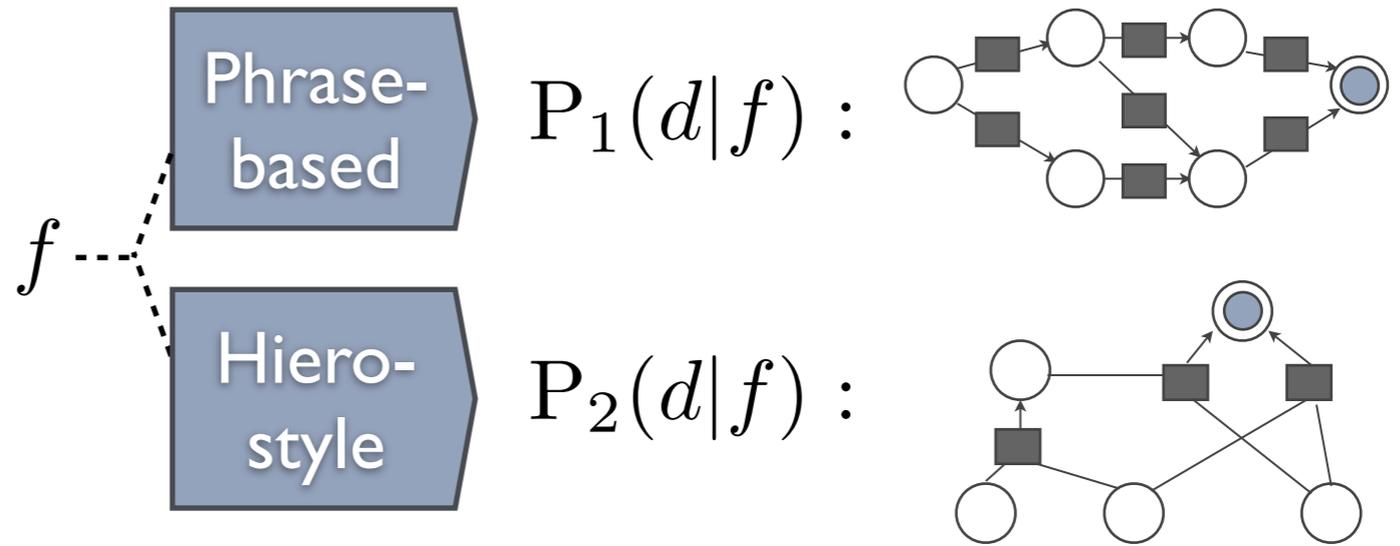
Extending Consensus Decoding to Multiple Models

I. Build posterior forests



Extending Consensus Decoding to Multiple Models

I. Build posterior forests



N-gram statistics

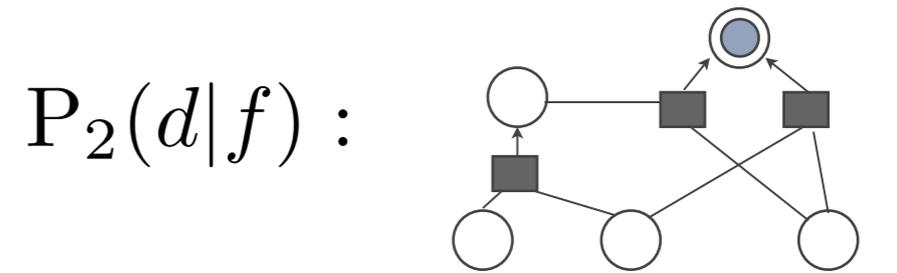
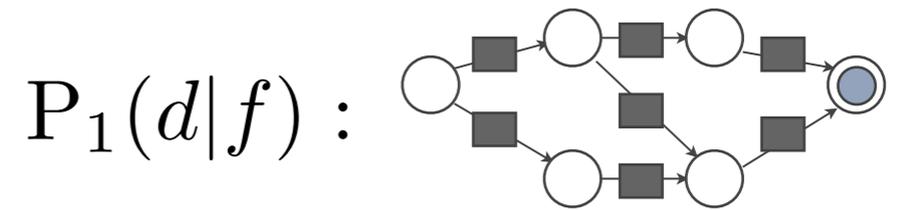
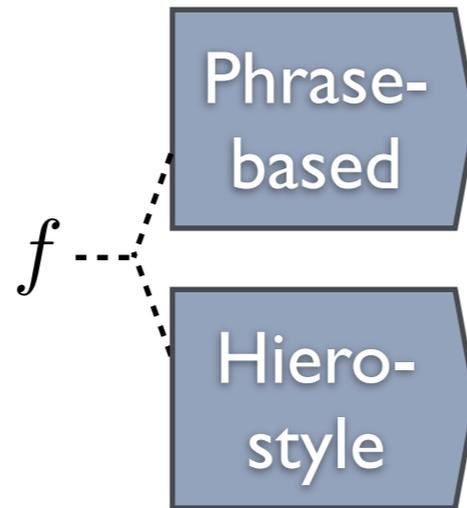
Length

Base model

$$C(d) = \sum_{n=1}^4 w^{(n)} \sum_{g \in n\text{-grams}} c(g, d) \cdot P(g|f) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} \theta \cdot \phi(d)$$

Extending Consensus Decoding to Multiple Models

I. Build posterior forests



N-gram statistics

$$C(d) = \sum_{n=1}^4 w^{(n)} v^{(n)}(d)$$

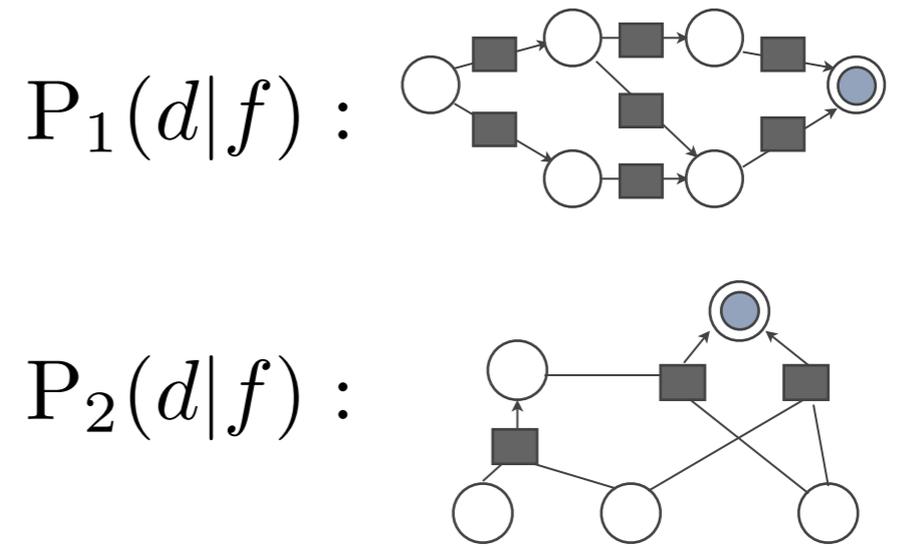
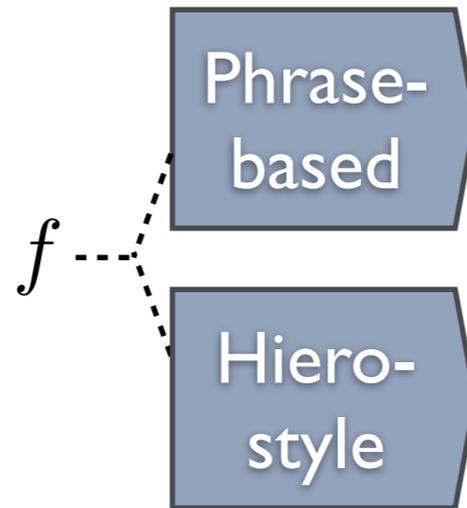
Length

$$+ w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

Base model

Extending Consensus Decoding to Multiple Models

I. Build posterior forests



Sum over models

N-gram statistics

Length

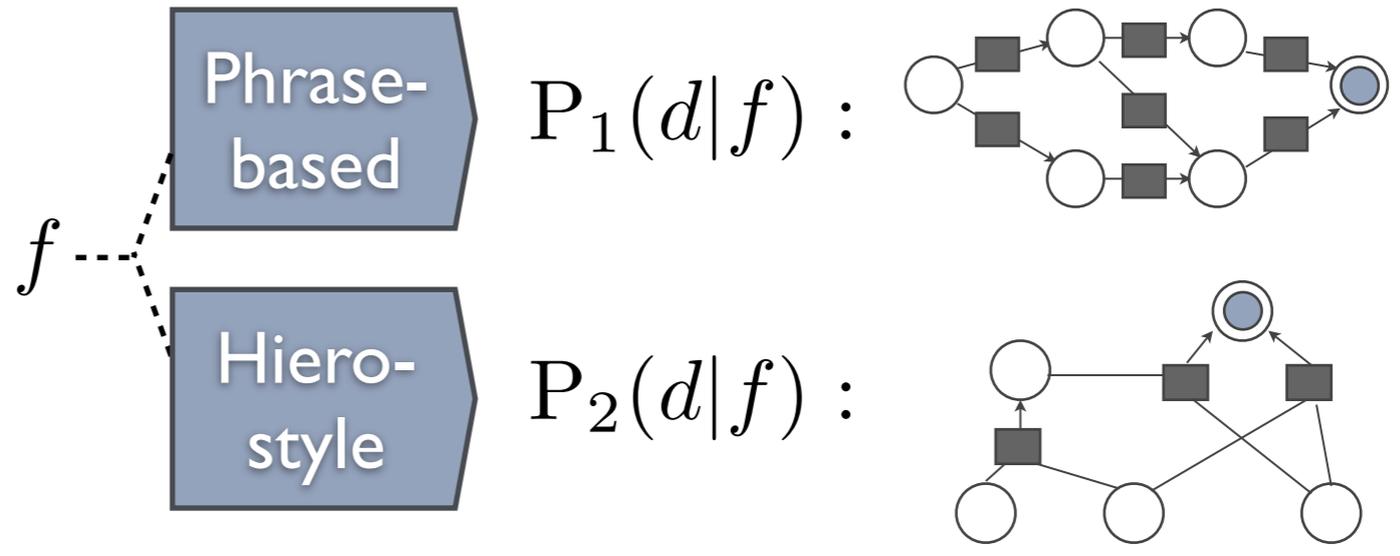
Base model

$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d)$$

$$+ w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

Extending Consensus Decoding to Multiple Models

1. Build posterior forests
2. Compute n-gram statistics from forests



Sum over models

N-gram statistics

$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d)$$

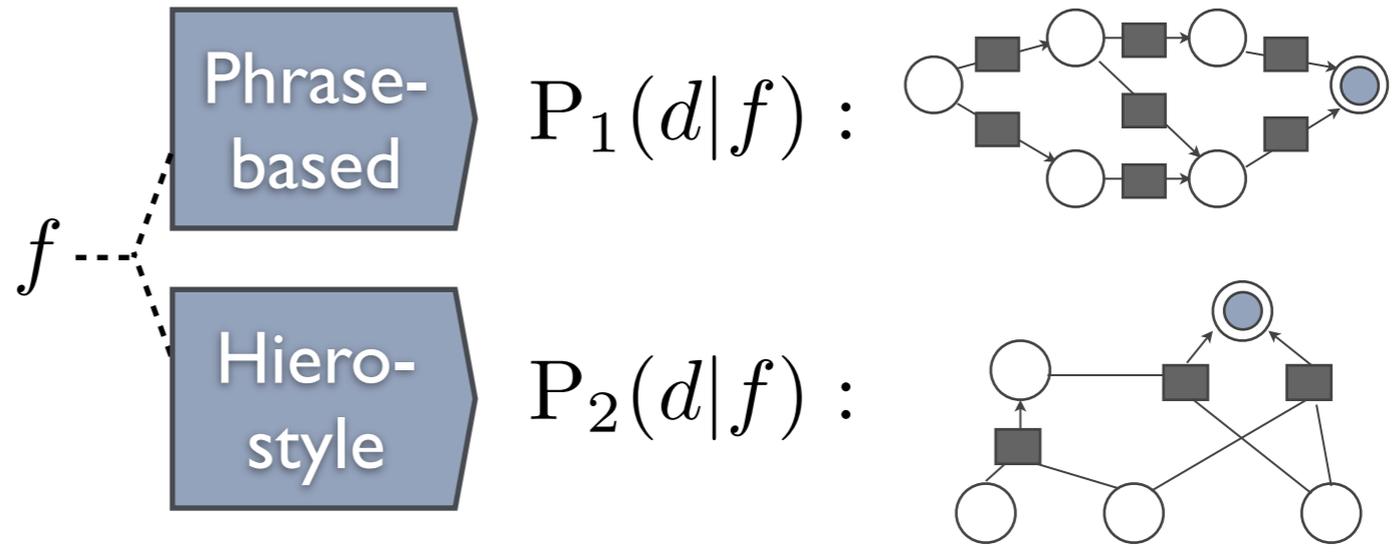
Length

Base model

$$+ w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

Extending Consensus Decoding to Multiple Models

1. Build posterior forests
2. Compute n-gram statistics from forests
3. Union all forests



Sum over models

N-gram statistics

$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d)$$

Length

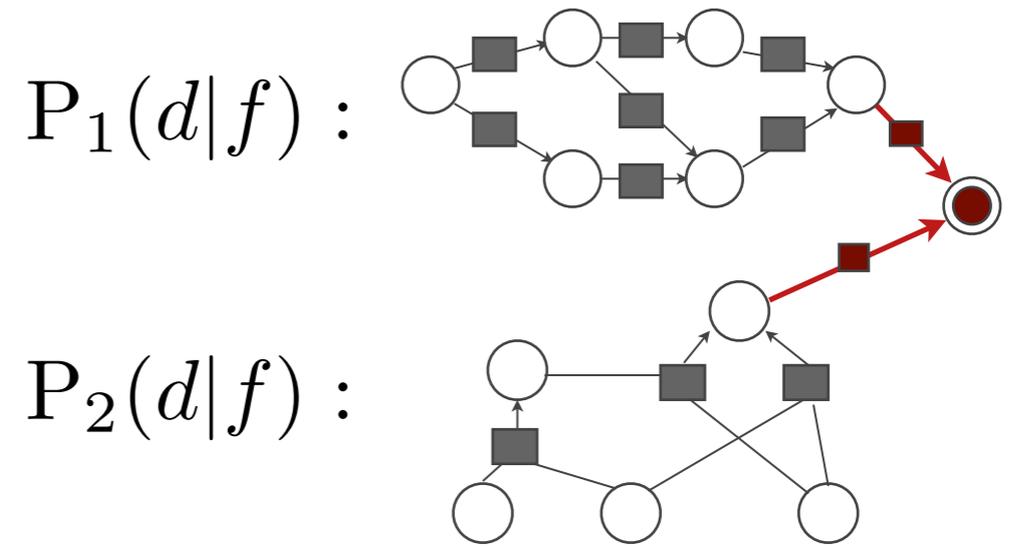
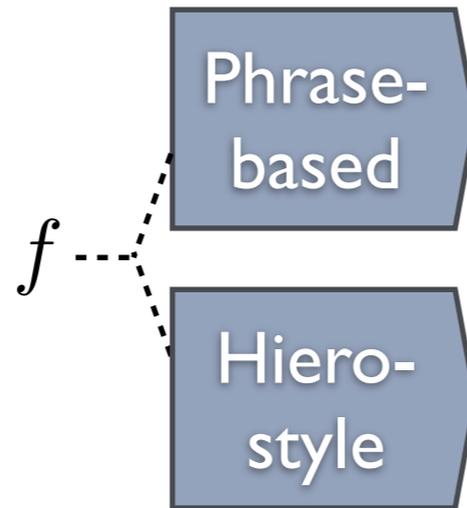
Base model

$$+ w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

Extending Consensus Decoding to Multiple Models

Google

1. Build posterior forests
2. Compute n-gram statistics from forests
3. Union all forests



Sum over models

N-gram statistics

$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d)$$

Length

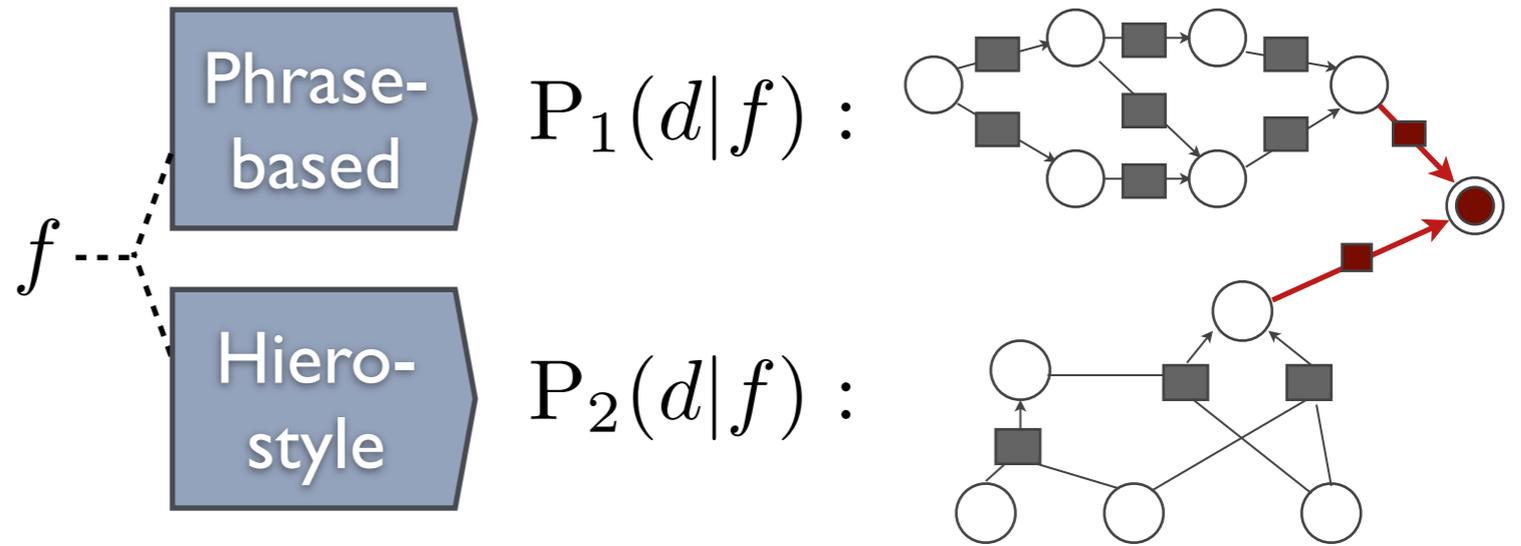
Base model

$$+ w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

Extending Consensus Decoding to Multiple Models

Google

1. Build posterior forests
2. Compute n-gram statistics from forests
3. Union all forests

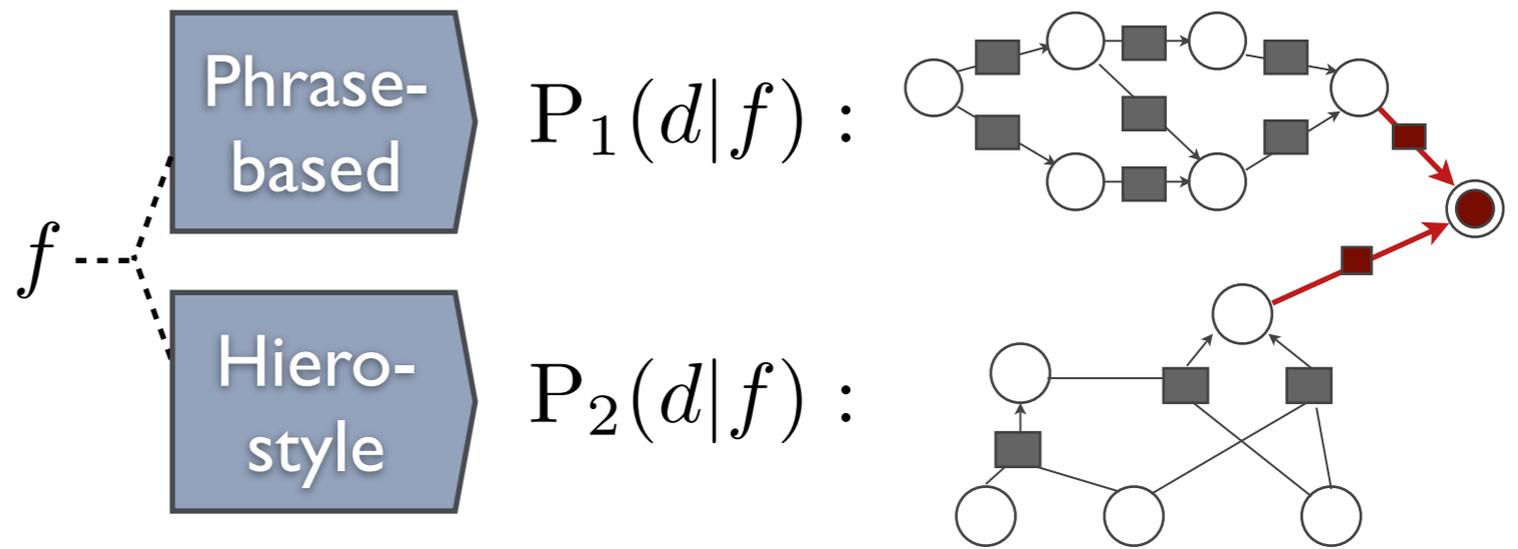


Sum over models N-gram statistics Model choice Length Base model

$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d) + w_i^{(\alpha)} \alpha_i(d) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

Extending Consensus Decoding to Multiple Models

1. Build posterior forests
2. Compute n-gram statistics from forests
3. Union all forests



$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d) + w_i^{(\alpha)} \alpha_i(d) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

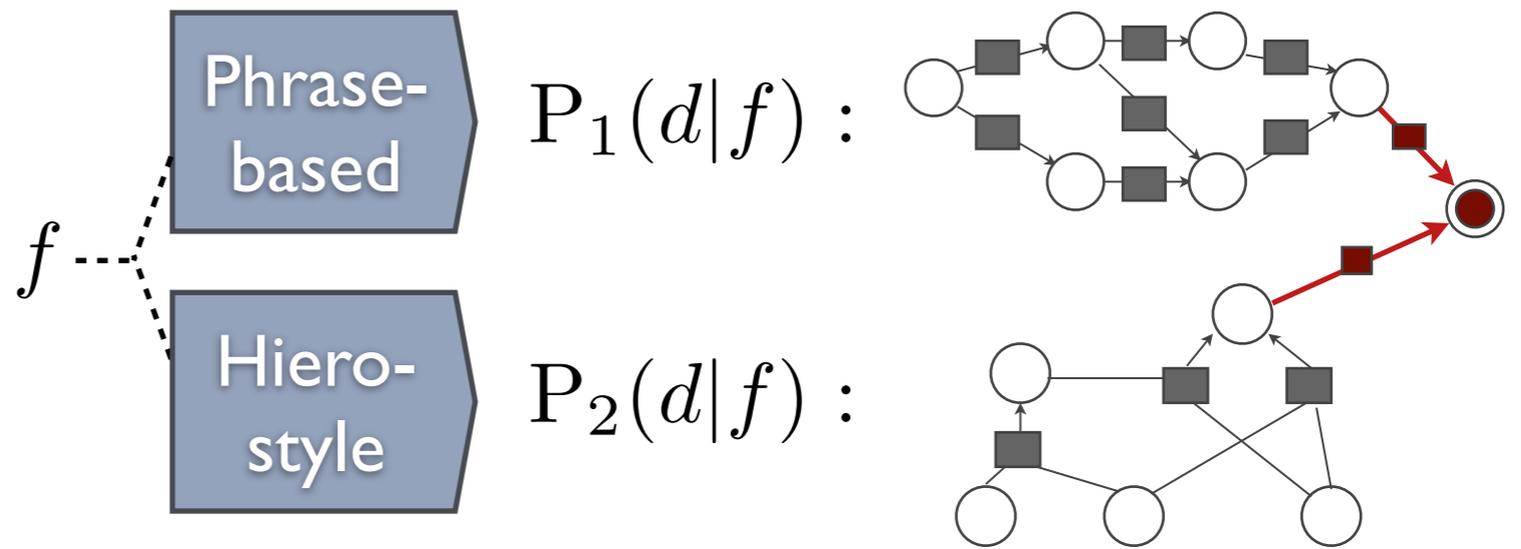
Sum over models N-gram statistics Model choice Length Base model

Model choice: Indicator feature for the system that originally generated d

Base model: Model score under the system that generated d

Extending Consensus Decoding to Multiple Models

1. Build posterior forests
2. Compute n-gram statistics from forests
3. Union all forests
4. Optimize multi-model consensus objective



$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d) + w_i^{(\alpha)} \alpha_i(d) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

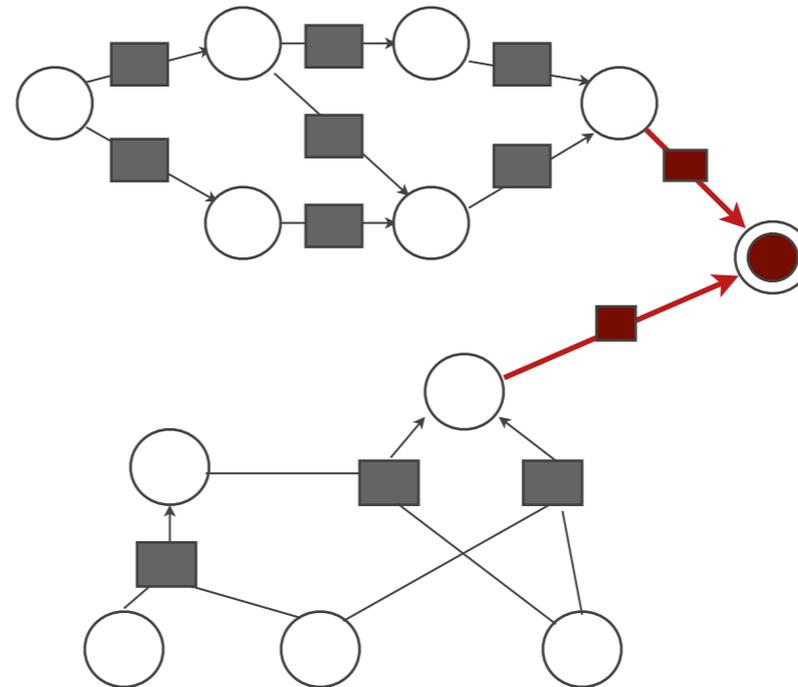
Sum over models N-gram statistics Model choice Length Base model

Model choice: Indicator feature for the system that originally generated d

Base model: Model score under the system that generated d

Extending Consensus Decoding to Multiple Models

1. Build posterior forests
2. Compute n-gram statistics from forests
3. Union all forests
4. Optimize multi-model consensus objective



Sum over models

N-gram statistics

Model choice

Length

Base model

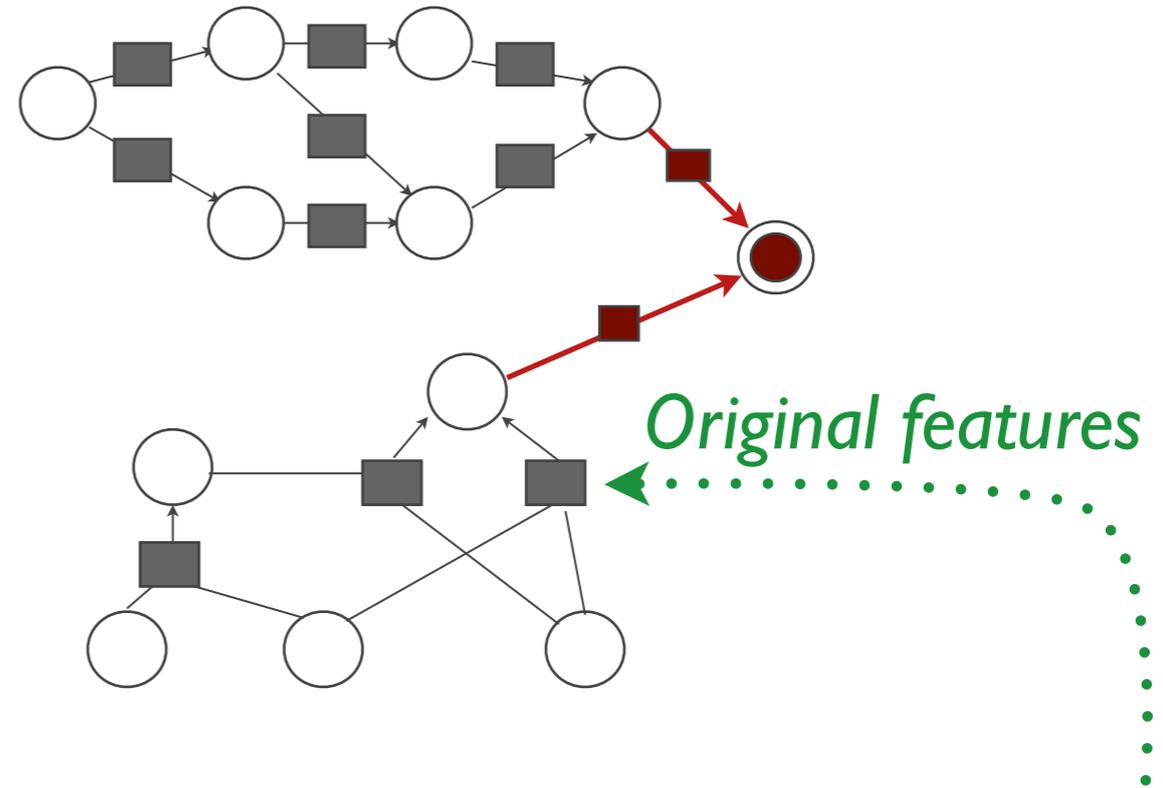
$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d) + w_i^{(\alpha)} \alpha_i(d) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

Model choice: Indicator feature for the system that originally generated d

Base model: Model score under the system that generated d

Extending Consensus Decoding to Multiple Models

1. Build posterior forests
2. Compute n-gram statistics from forests
3. Union all forests
4. Optimize multi-model consensus objective



$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d) + w_i^{(\alpha)} \alpha_i(d) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

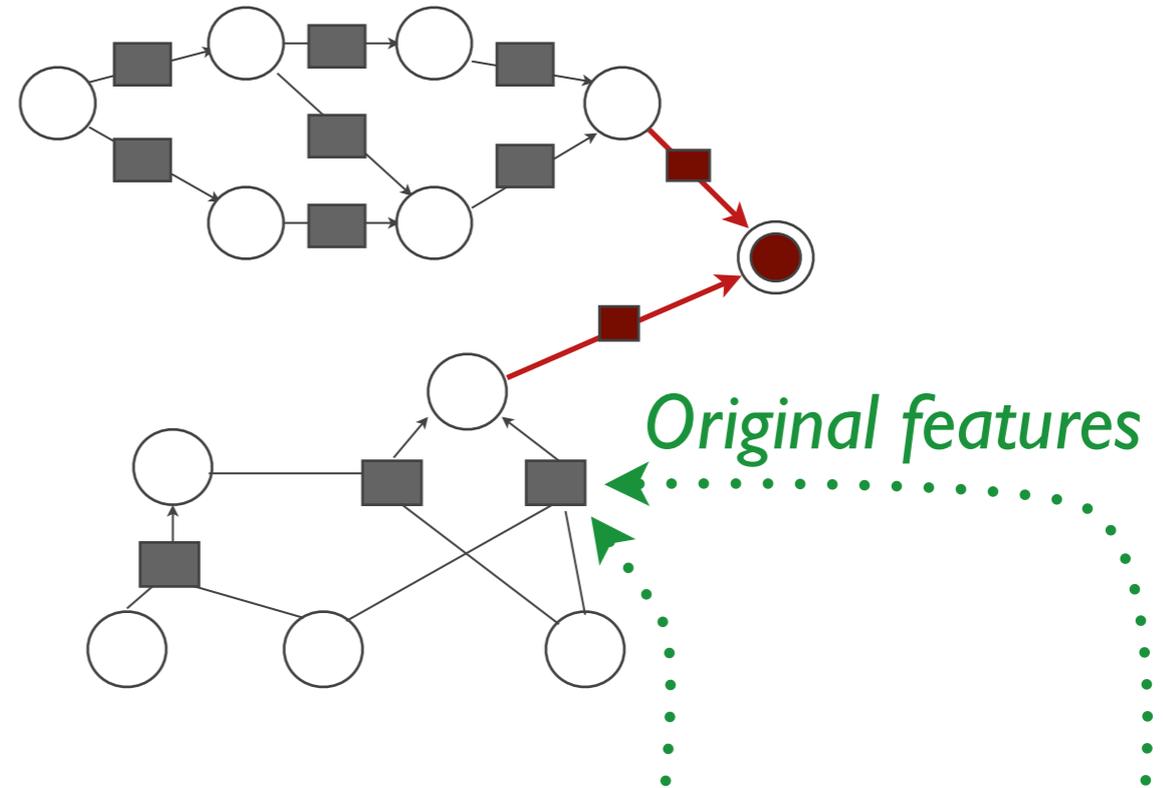
Sum over models
N-gram statistics
Model choice
Length
Base model

Model choice: Indicator feature for the system that originally generated d

Base model: Model score under the system that generated d

Extending Consensus Decoding to Multiple Models

1. Build posterior forests
2. Compute n-gram statistics from forests
3. Union all forests
4. Optimize multi-model consensus objective



Sum over models N-gram statistics Model choice Length Base model

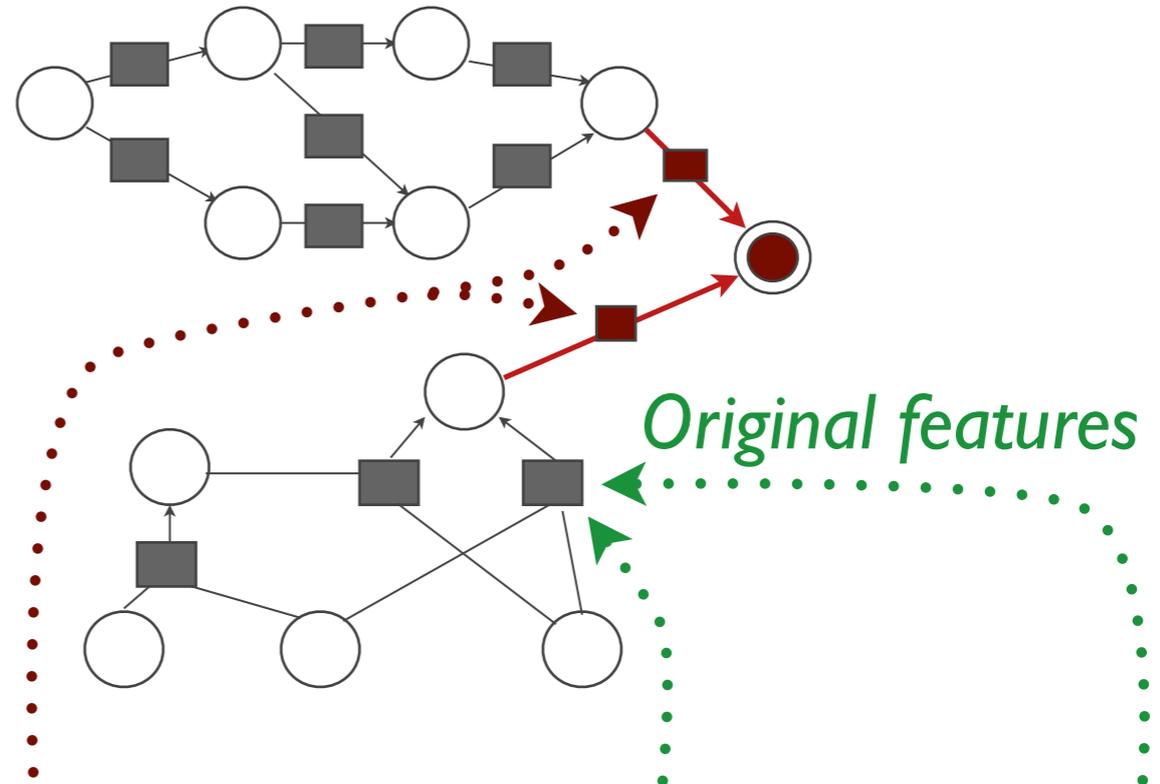
$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d) + w_i^{(\alpha)} \alpha_i(d) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

Model choice: Indicator feature for the system that originally generated d

Base model: Model score under the system that generated d

Extending Consensus Decoding to Multiple Models

1. Build posterior forests
2. Compute n-gram statistics from forests
3. Union all forests
4. Optimize multi-model consensus objective



Sum over models N-gram statistics Model choice Length Base model

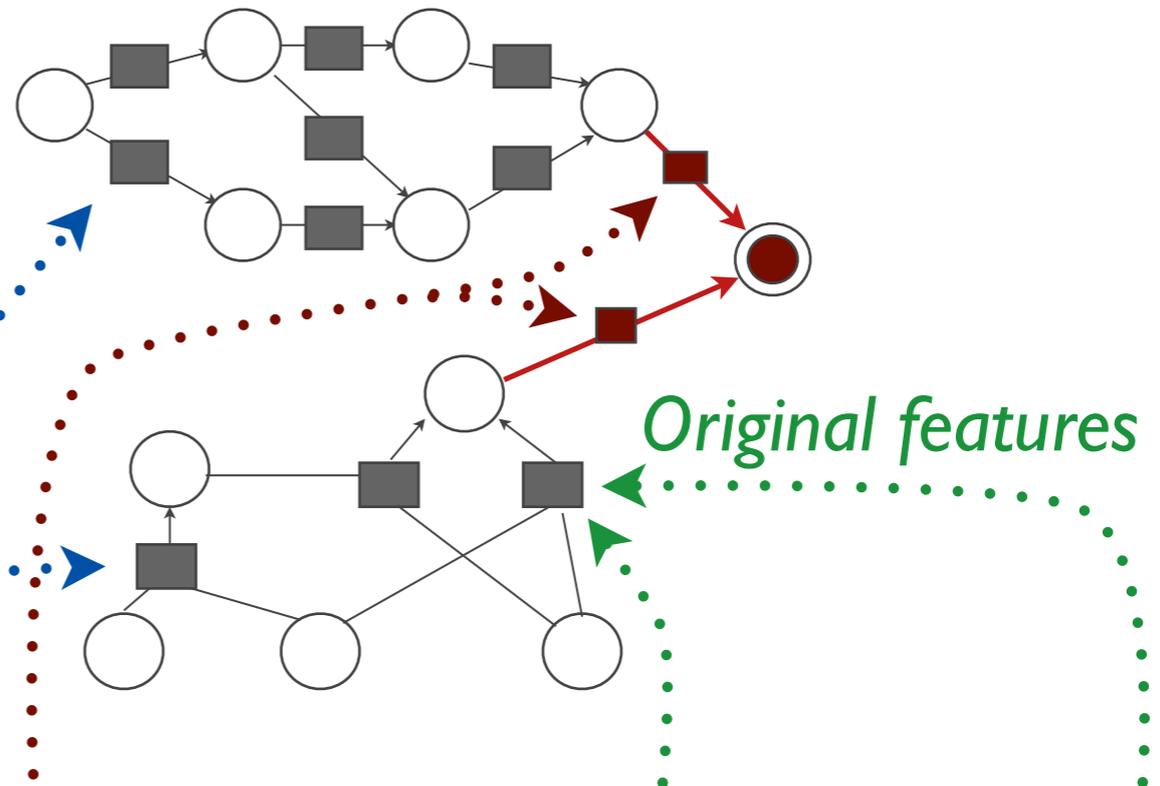
$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d) + w_i^{(\alpha)} \alpha_i(d) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

Model choice: Indicator feature for the system that originally generated d

Base model: Model score under the system that generated d

Extending Consensus Decoding to Multiple Models

1. Build posterior forests
2. Compute n-gram statistics from forests
3. Union all forests
4. Optimize multi-model consensus objective



Sum over models N-gram statistics Model choice Length Base model

$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d) + w_i^{(\alpha)} \alpha_i(d) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

Model choice: Indicator feature for the system that originally generated d

Base model: Model score under the system that generated d

Properties of Model Combination

Google

Sum over models N-gram statistics Model choice Length Base model

$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d) + w_i^{(\alpha)} \alpha_i(d) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

Properties of Model Combination

Google

Sum over models N-gram statistics Model choice Length Base model

$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d) + w_i^{(\alpha)} \alpha_i(d) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

- ▶ Reduces to consensus decoding when we have only one model

Properties of Model Combination

Google

Sum over models N-gram statistics Model choice Length Base model

$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d) + w_i^{(\alpha)} \alpha_i(d) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

- ▶ Reduces to consensus decoding when we have only one model
- ▶ A linear model: w can be tuned to maximize output performance

Properties of Model Combination

Google

Sum over models N-gram statistics Model choice Length Base model

$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d) + w_i^{(\alpha)} \alpha_i(d) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

- ▶ Reduces to consensus decoding when we have only one model
- ▶ A linear model: w can be tuned to maximize output performance
- ▶ No concept of a primary system

Properties of Model Combination

Google

Sum over models N-gram statistics Model choice Length Base model

$$C(d) = \sum_{i=1}^I \sum_{n=1}^4 w_i^{(n)} v_i^{(n)}(d) + w_i^{(\alpha)} \alpha_i(d) + w^{(\ell)} |\sigma_e(d)| + w^{(b)} b(d)$$

- ▶ Reduces to consensus decoding when we have only one model
- ▶ A linear model: w can be tuned to maximize output performance
- ▶ No concept of a primary system
- ▶ Every possible output was a derivation under *some original model*

Model Combination Experimental Results

Google

	Max-derivation
	Consensus

- ▶ Compared three in-house Google systems
- ▶ Constrained data track of 2008 NIST task
- ▶ Parameters tuned on NIST 2004 eval set

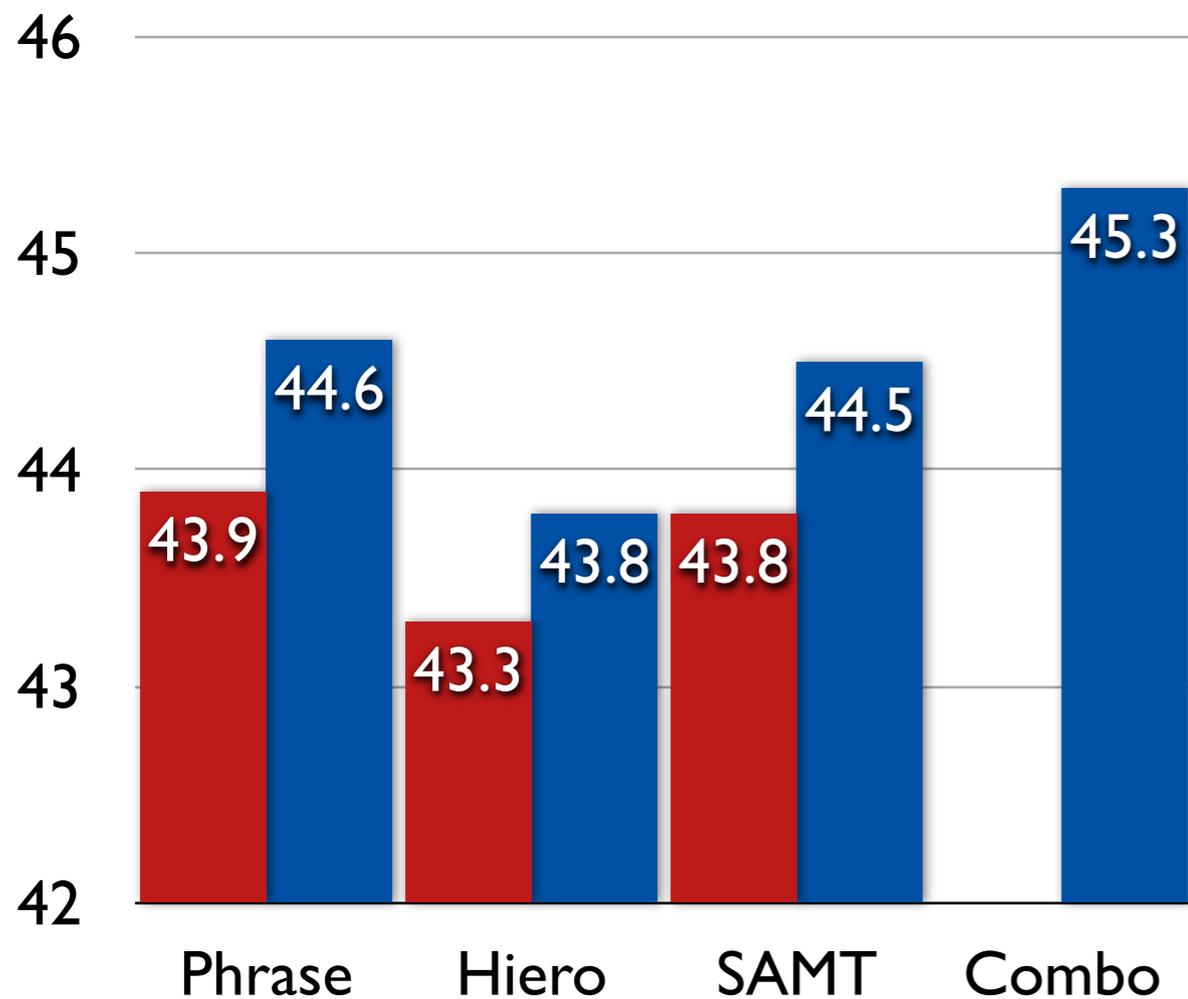
Model Combination Experimental Results

Google



- ▶ Compared three in-house Google systems
- ▶ Constrained data track of 2008 NIST task
- ▶ Parameters tuned on NIST 2004 eval set

Arabic-to-English BLEU



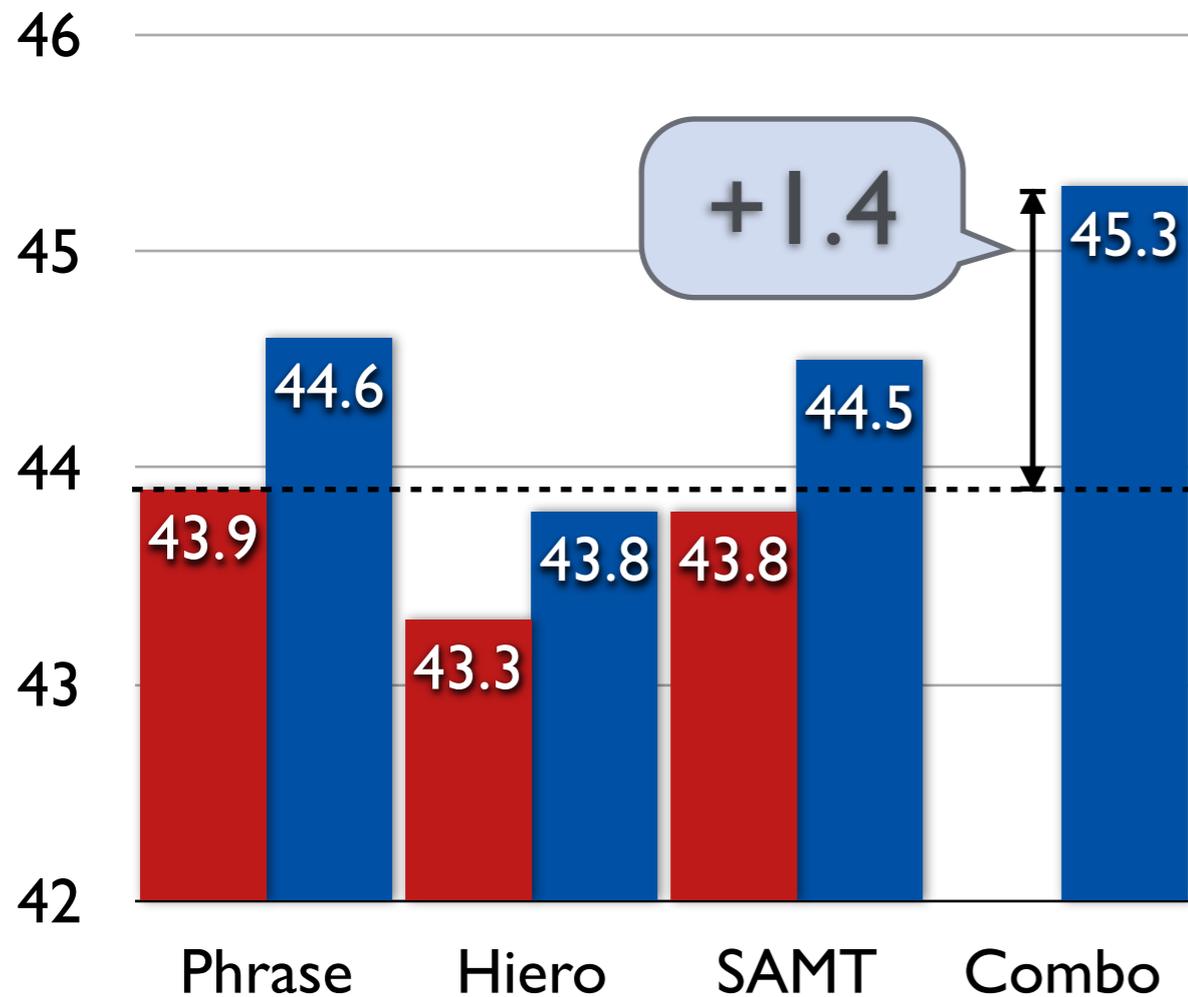
Model Combination Experimental Results

Google



- ▶ Compared three in-house Google systems
- ▶ Constrained data track of 2008 NIST task
- ▶ Parameters tuned on NIST 2004 eval set

Arabic-to-English BLEU



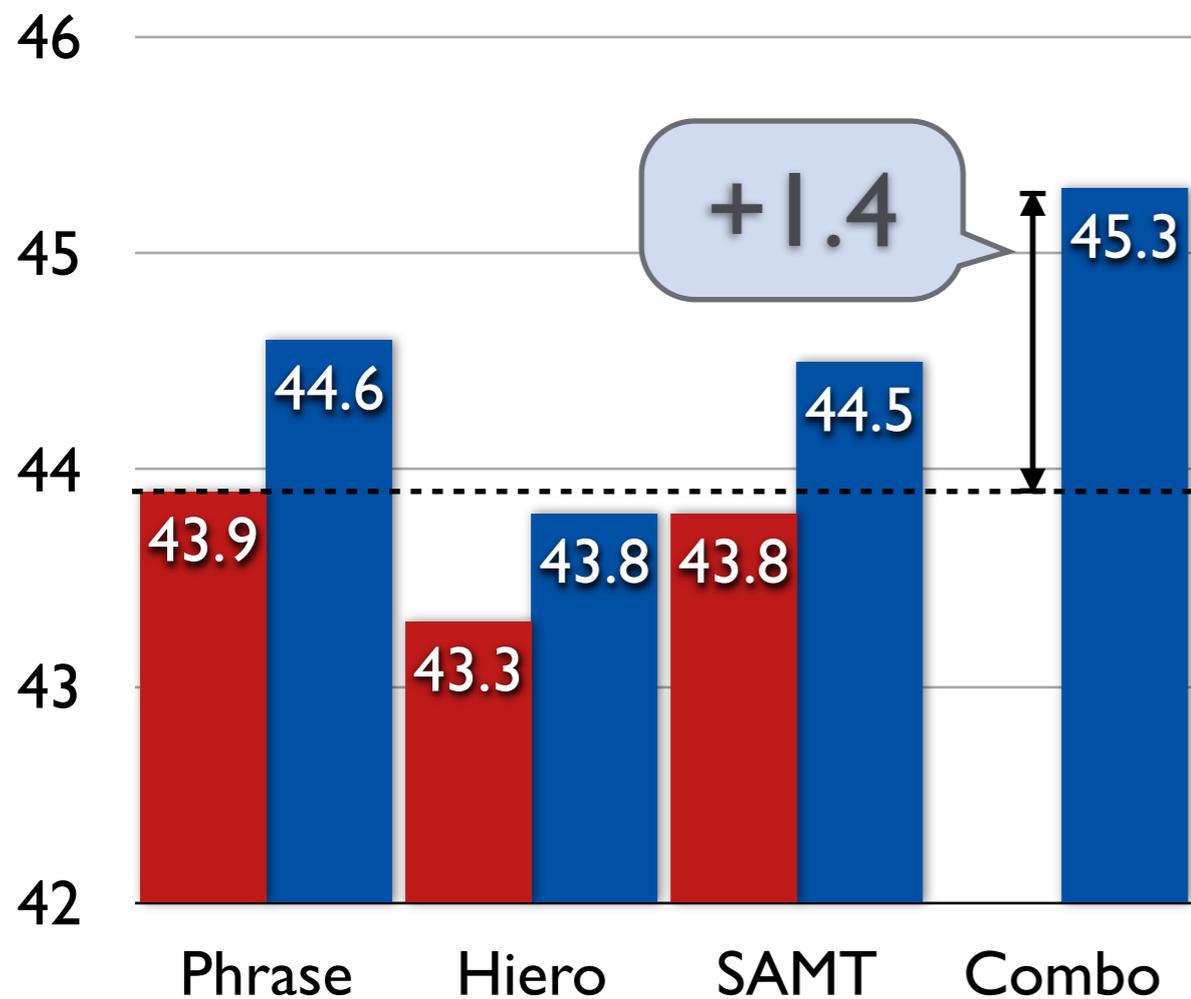
Model Combination Experimental Results

Google

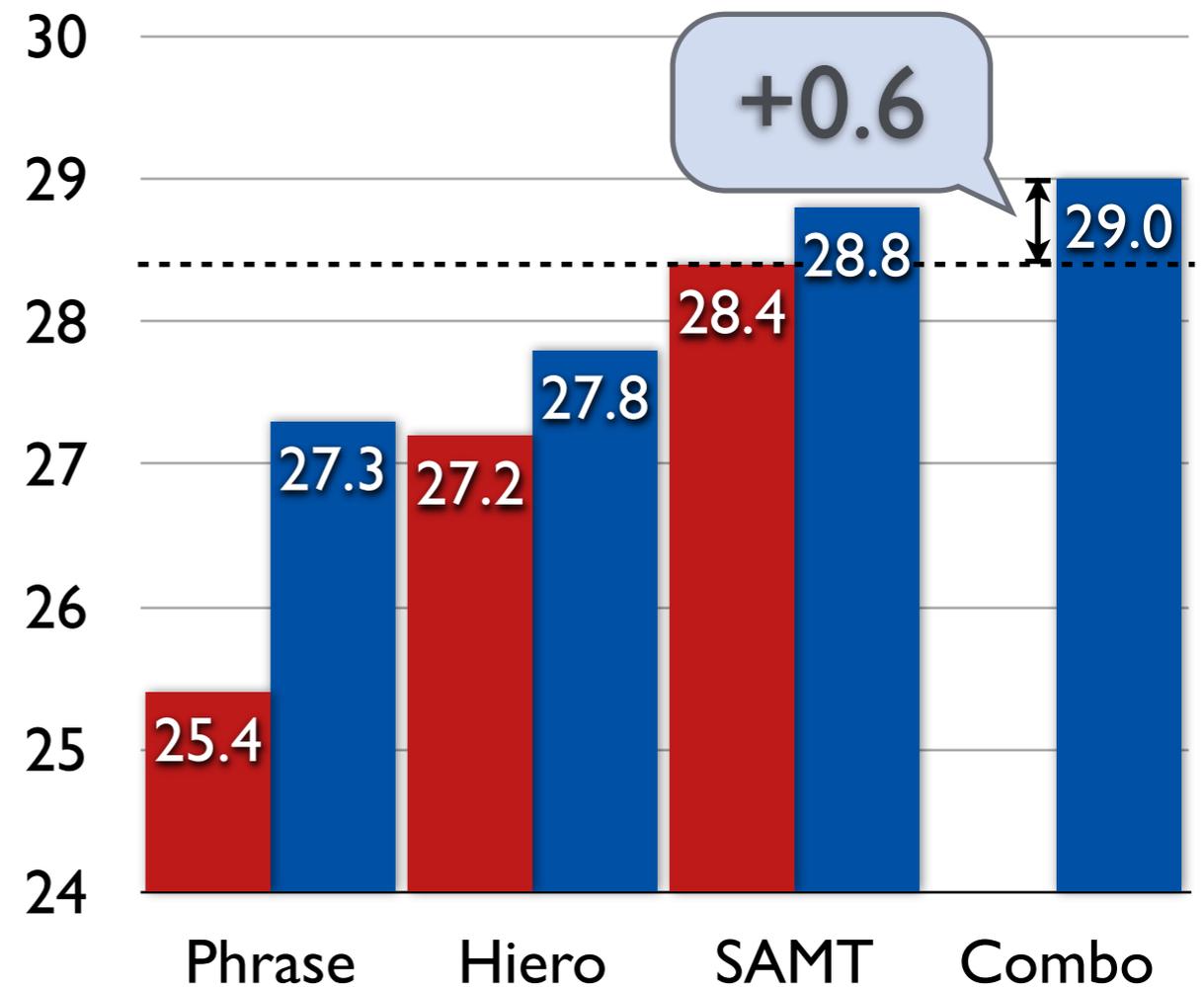


- ▶ Compared three in-house Google systems
- ▶ Constrained data track of 2008 NIST task
- ▶ Parameters tuned on NIST 2004 eval set

Arabic-to-English BLEU



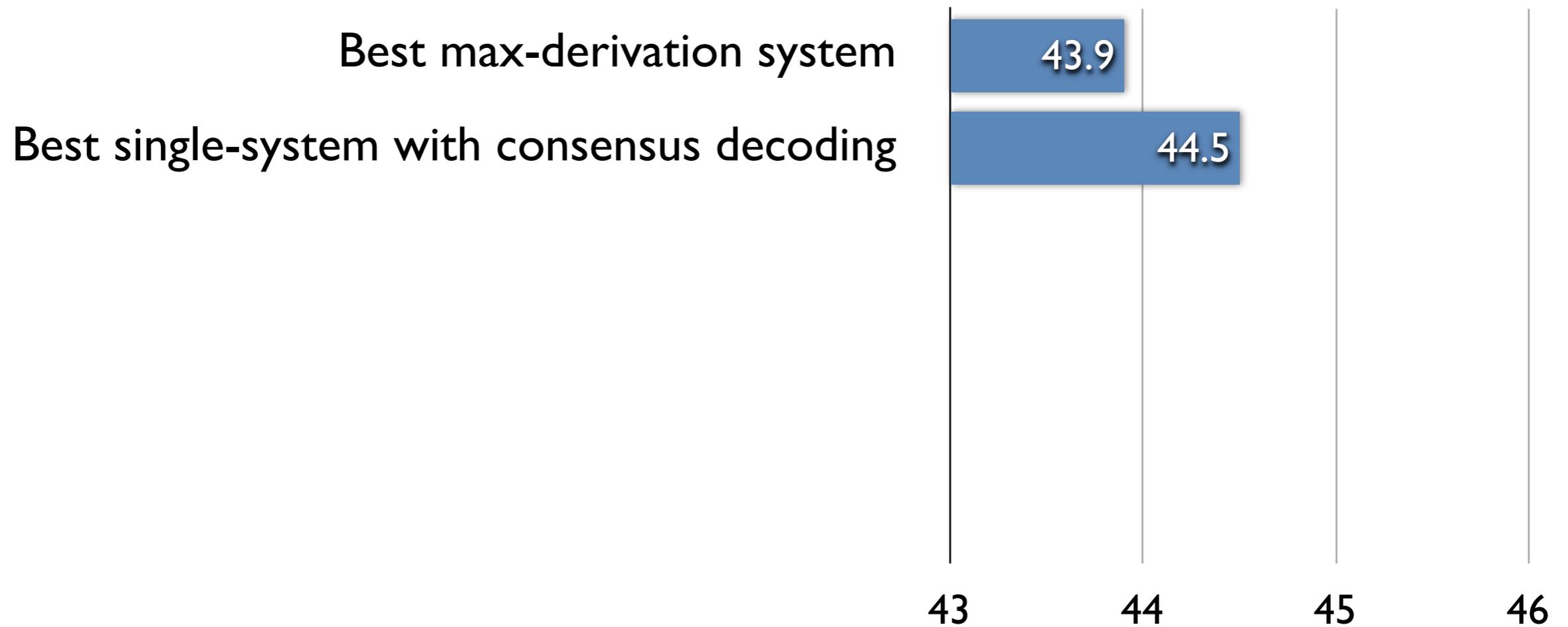
Chinese-to-English BLEU



Sources of Improvement: Arabic-to-English

Google

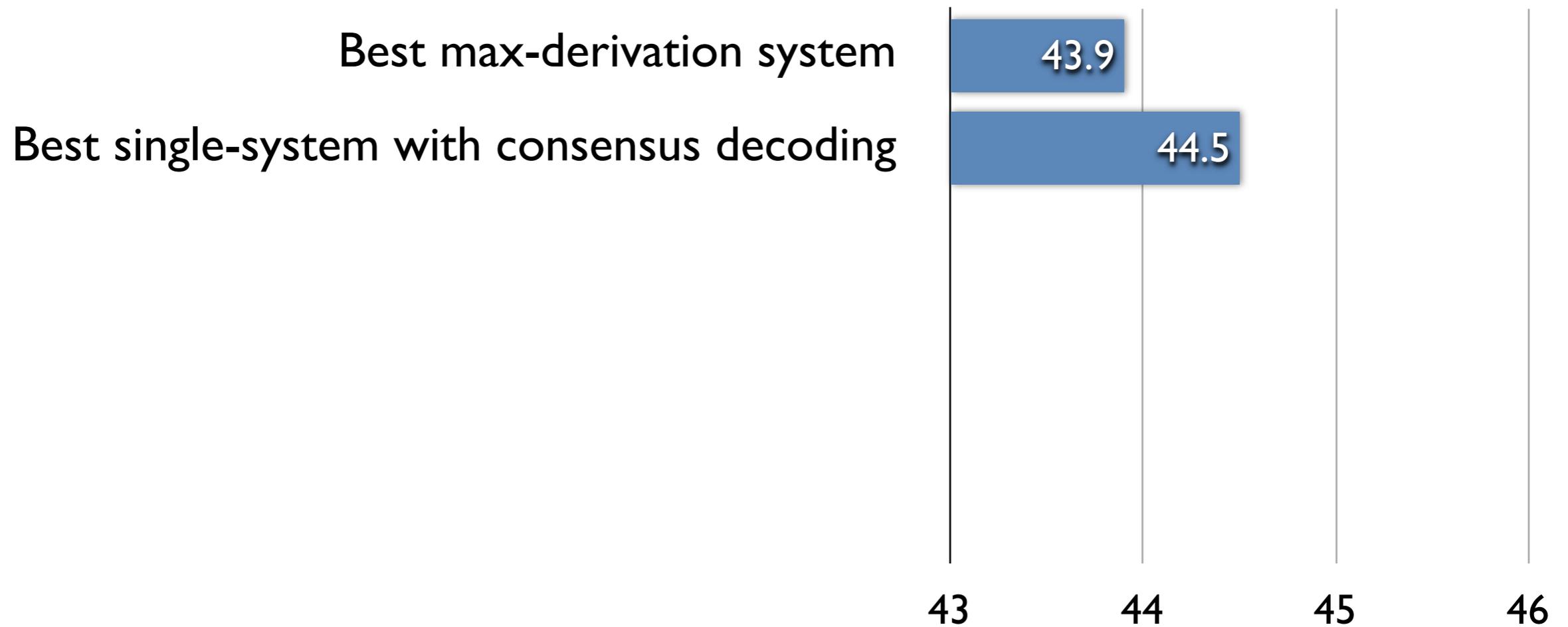
Sources of Improvement: Arabic-to-English



Sources of Improvement: Arabic-to-English



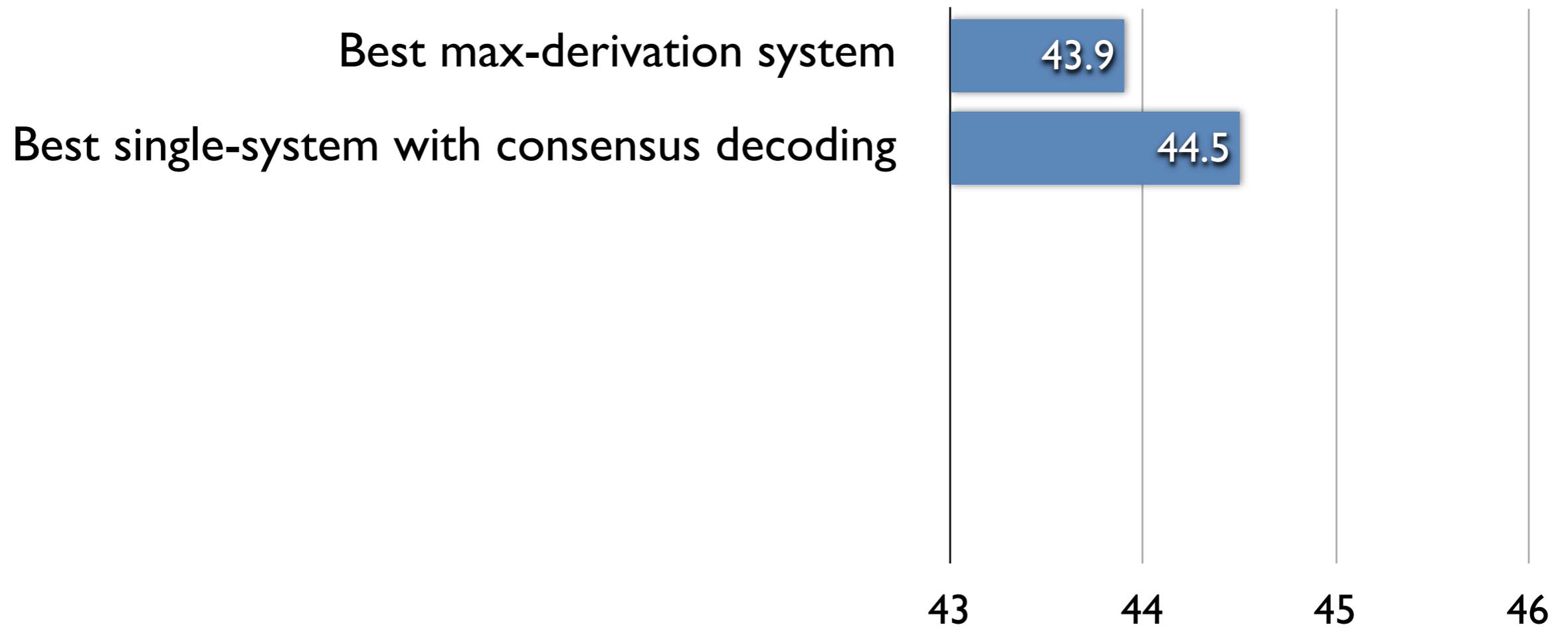
- ▶ Do we need model choice features?



Sources of Improvement: Arabic-to-English

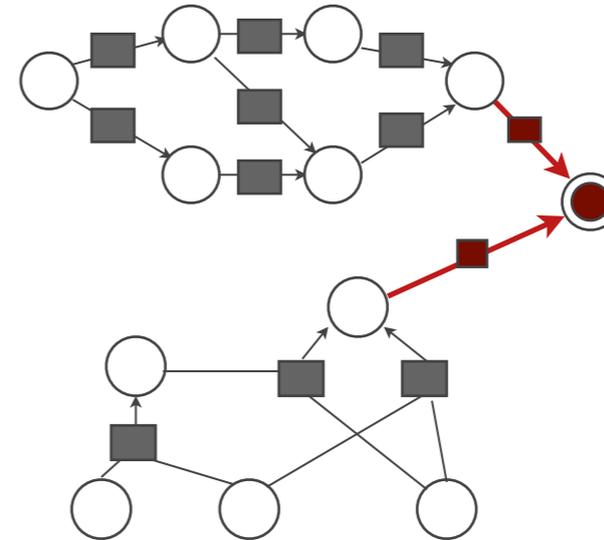
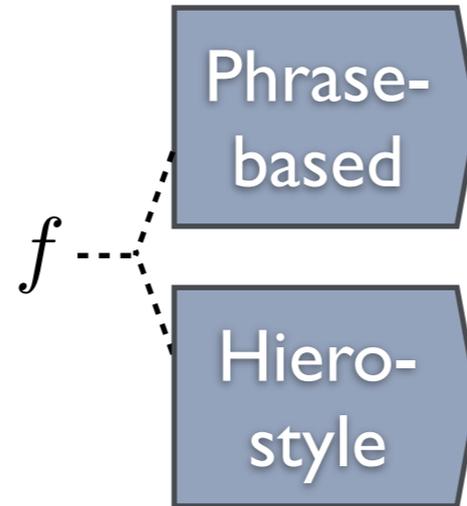


- ▶ Do we need model choice features?
- ▶ Should n-gram statistics be model specific?

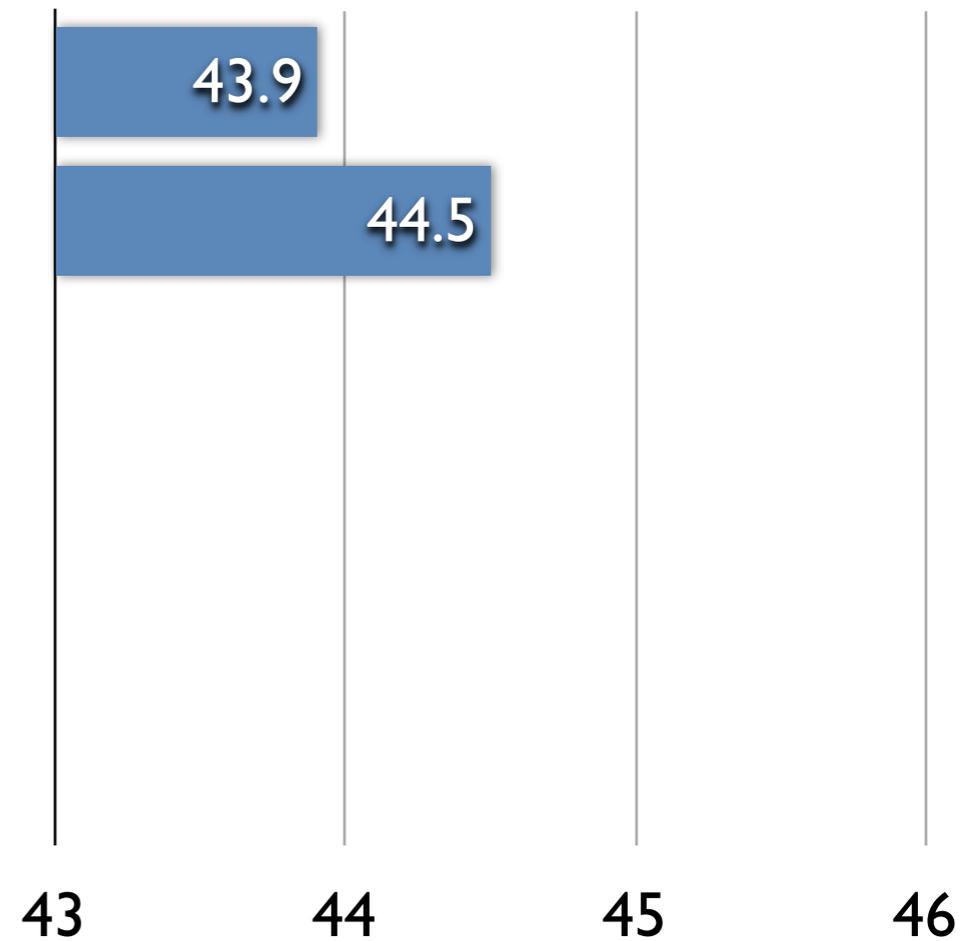


Sources of Improvement: Arabic-to-English

- ▶ Do we need model choice features?
- ▶ Should n-gram statistics be model specific?

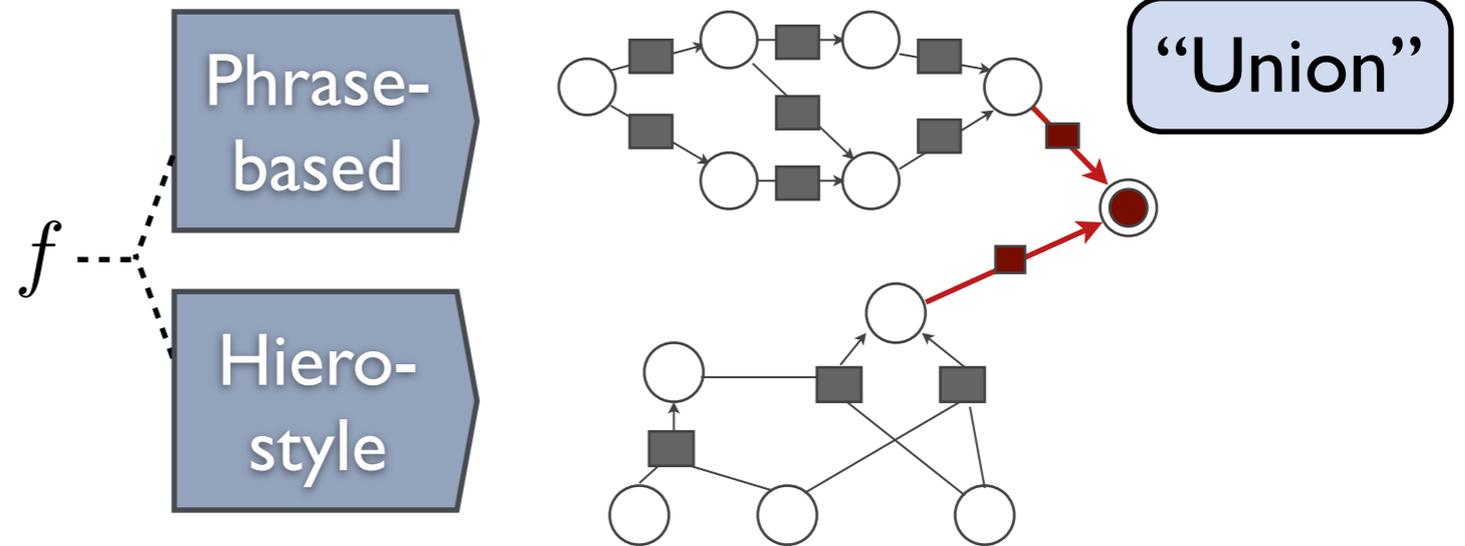


Best max-derivation system
Best single-system with consensus decoding

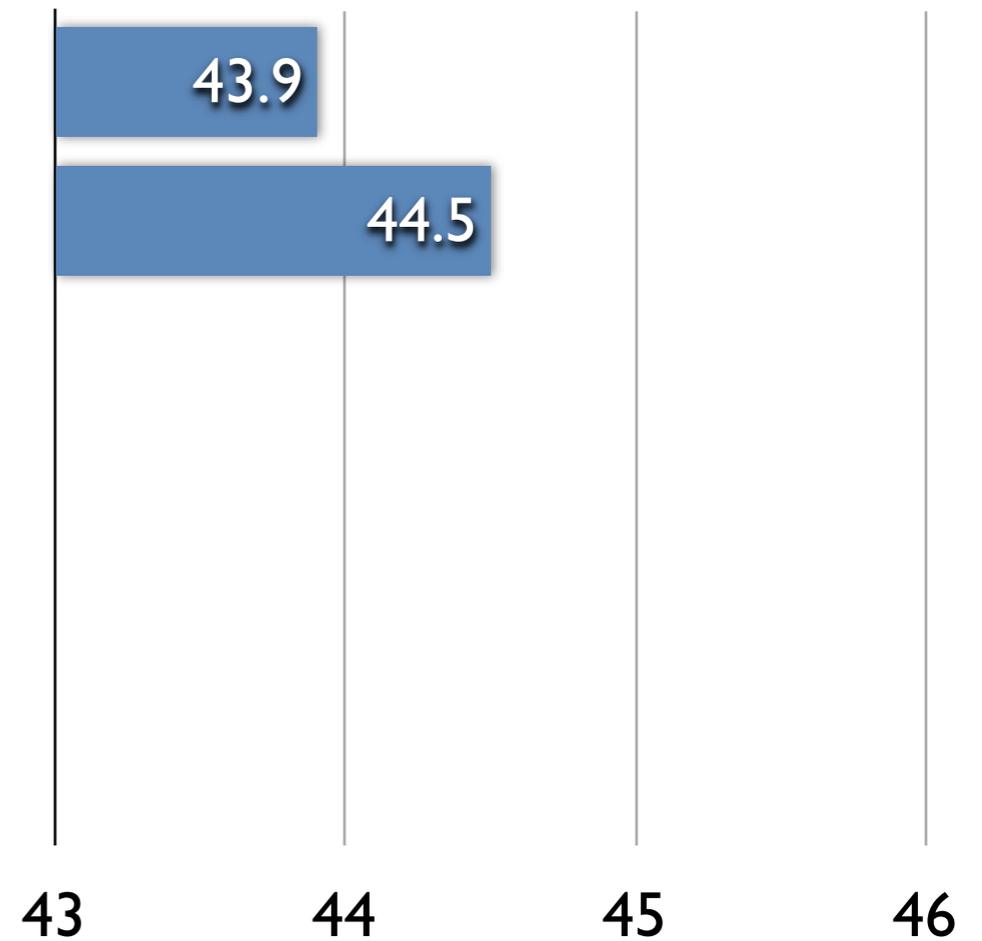


Sources of Improvement: Arabic-to-English

- ▶ Do we need model choice features?
- ▶ Should n-gram statistics be model specific?

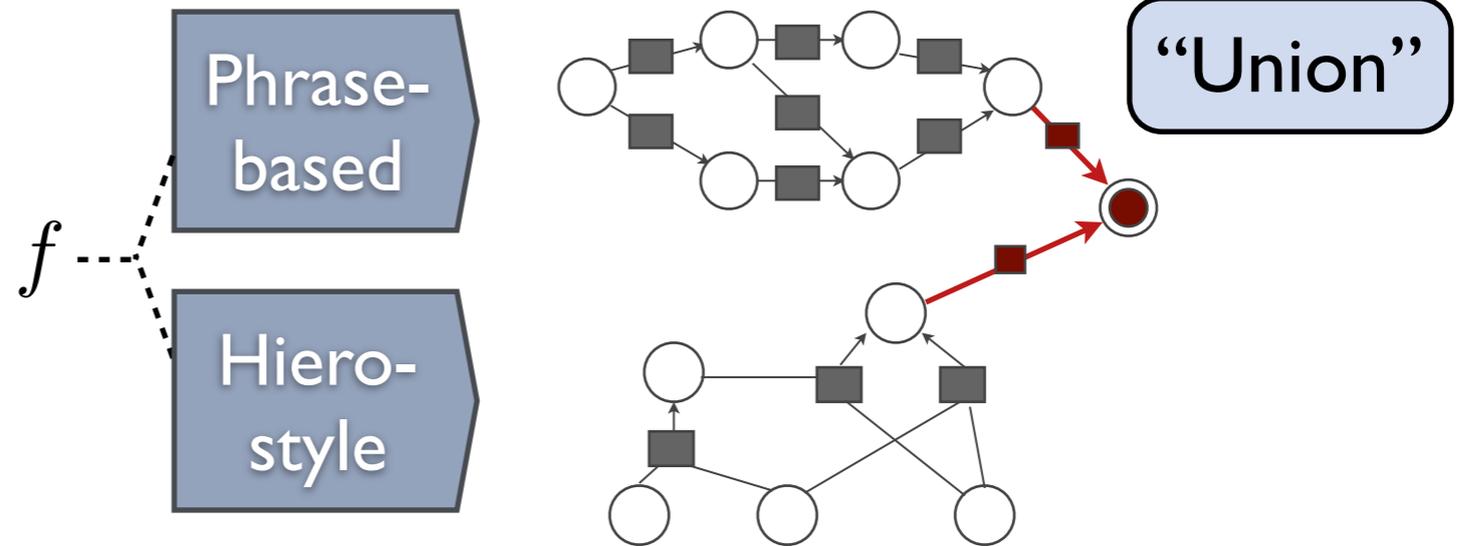


Best max-derivation system
Best single-system with consensus decoding

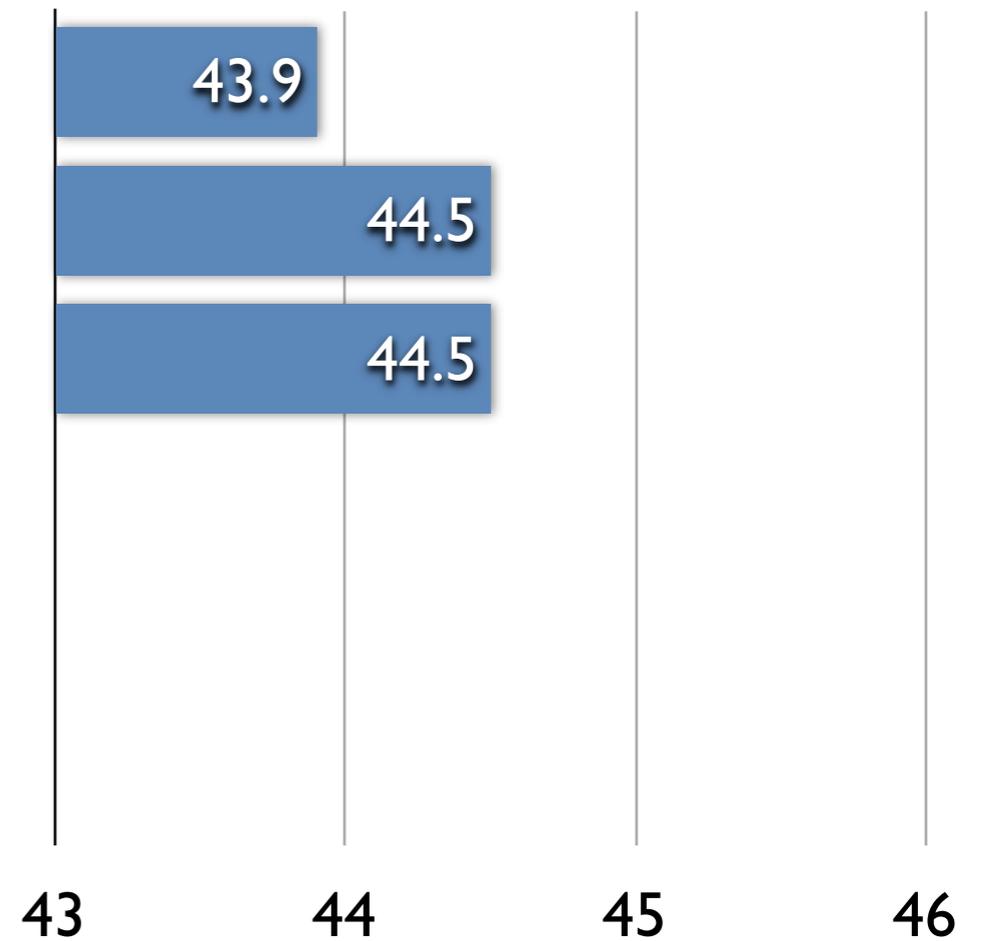


Sources of Improvement: Arabic-to-English

- ▶ Do we need model choice features?
- ▶ Should n-gram statistics be model specific?

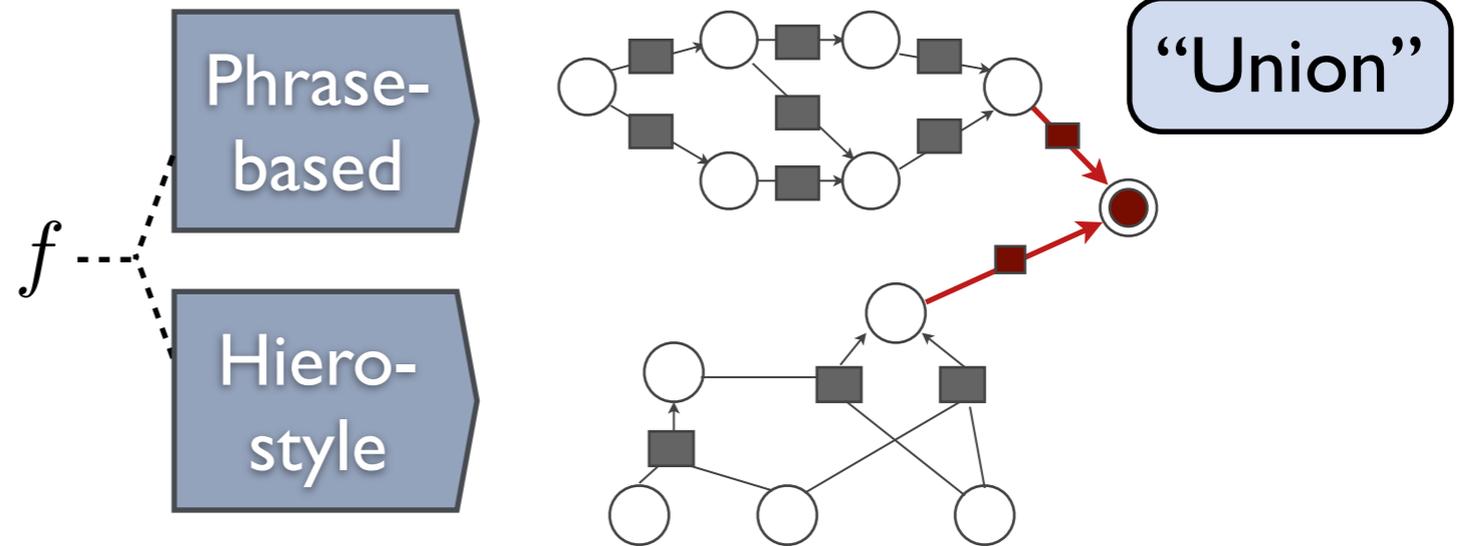


Best max-derivation system
Best single-system with consensus decoding
Union + consensus decoding

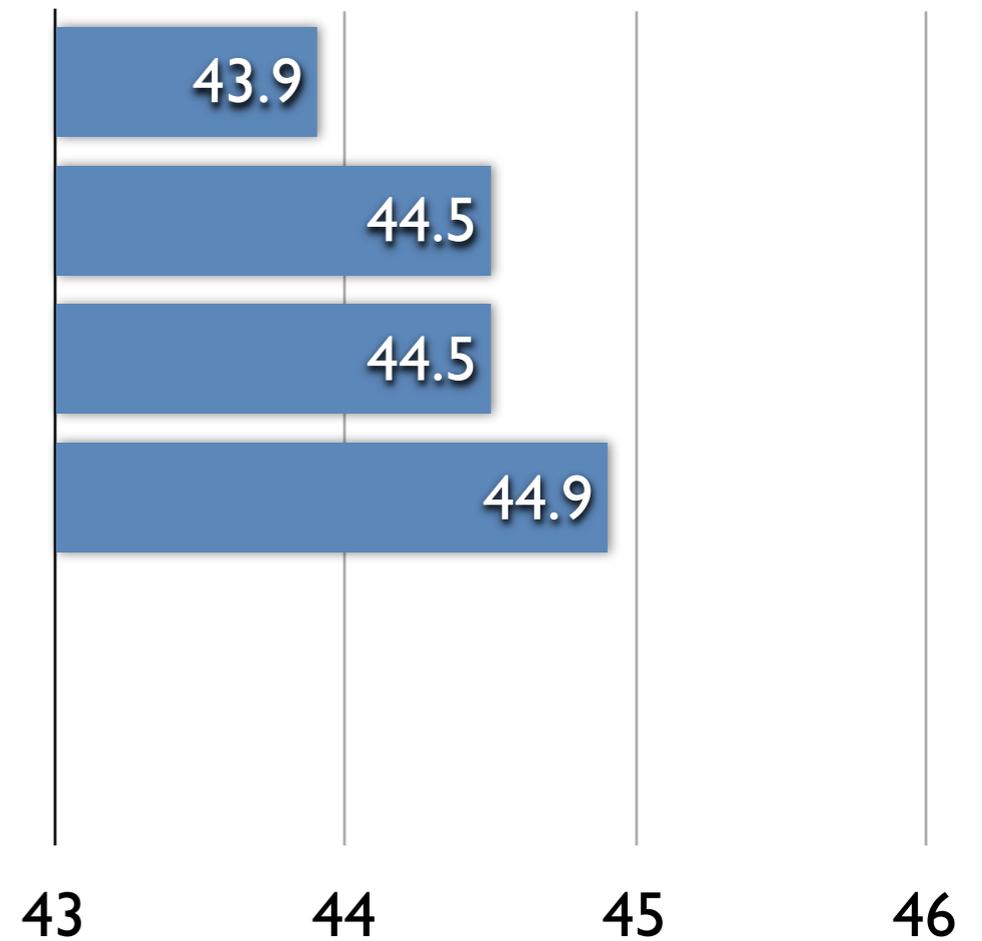


Sources of Improvement: Arabic-to-English

- ▶ Do we need model choice features?
- ▶ Should n-gram statistics be model specific?



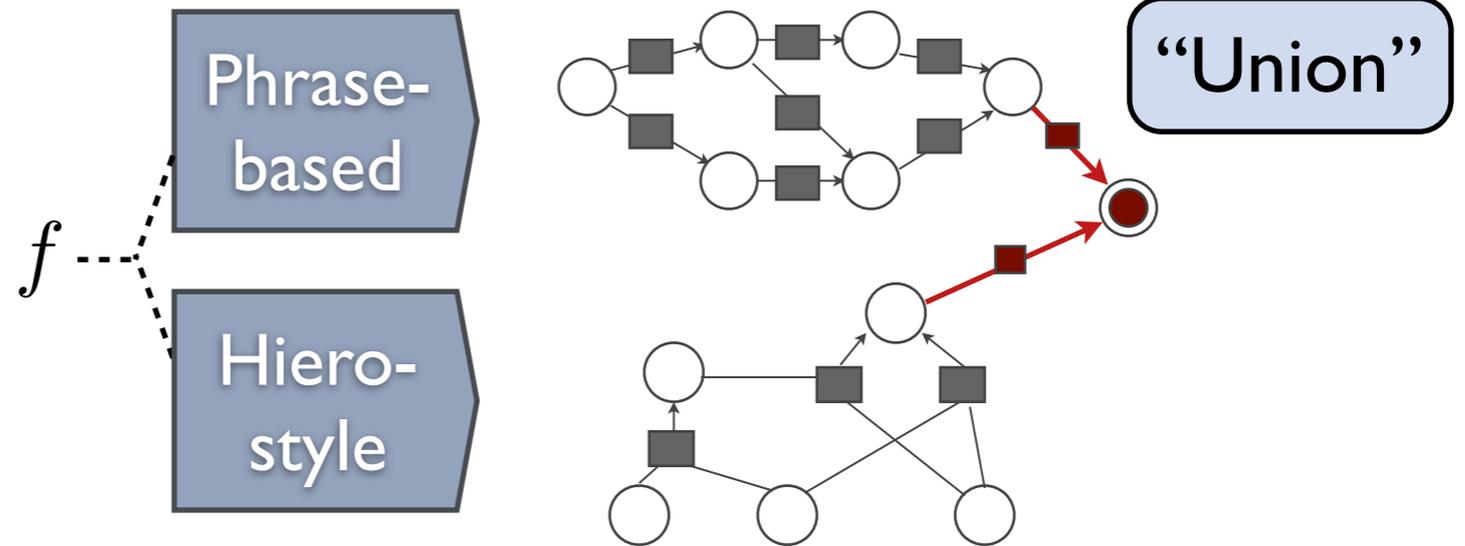
Best max-derivation system
Best single-system with consensus decoding
Union + consensus decoding
Union + consensus + model choice features



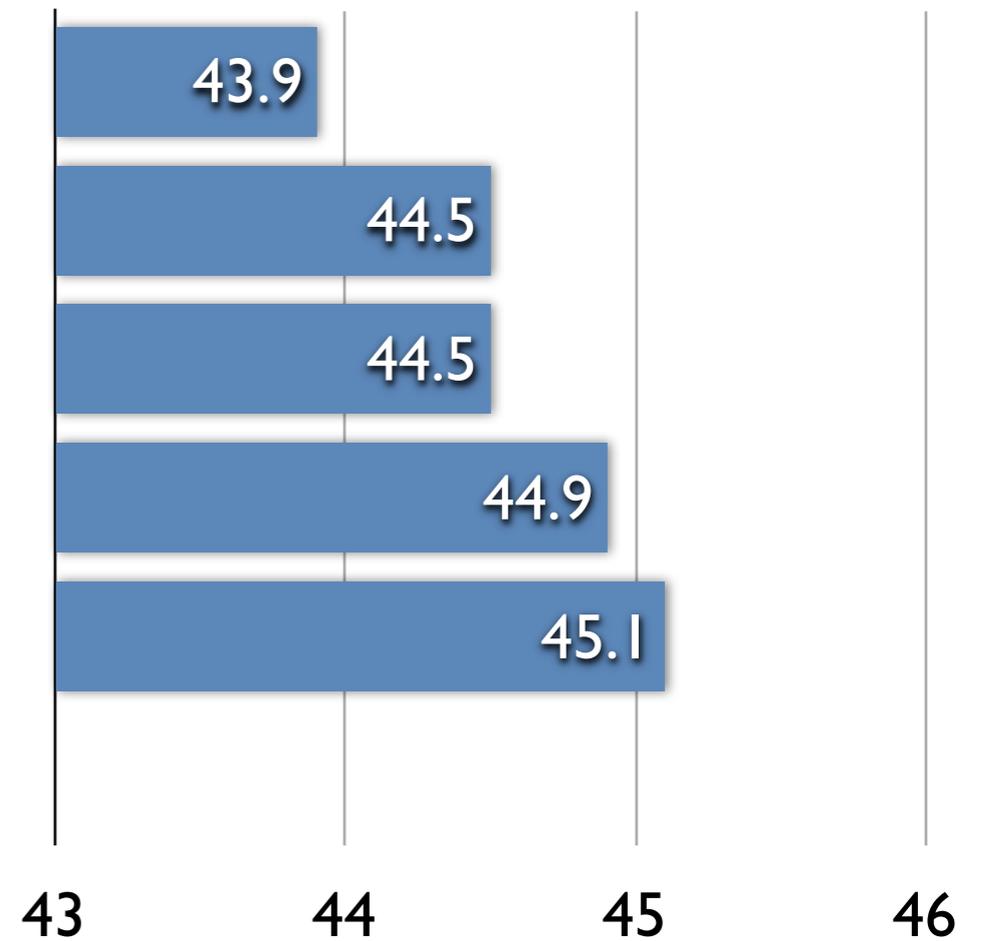
Sources of Improvement: Arabic-to-English

Google

- ▶ Do we need model choice features?
- ▶ Should n-gram statistics be model specific?



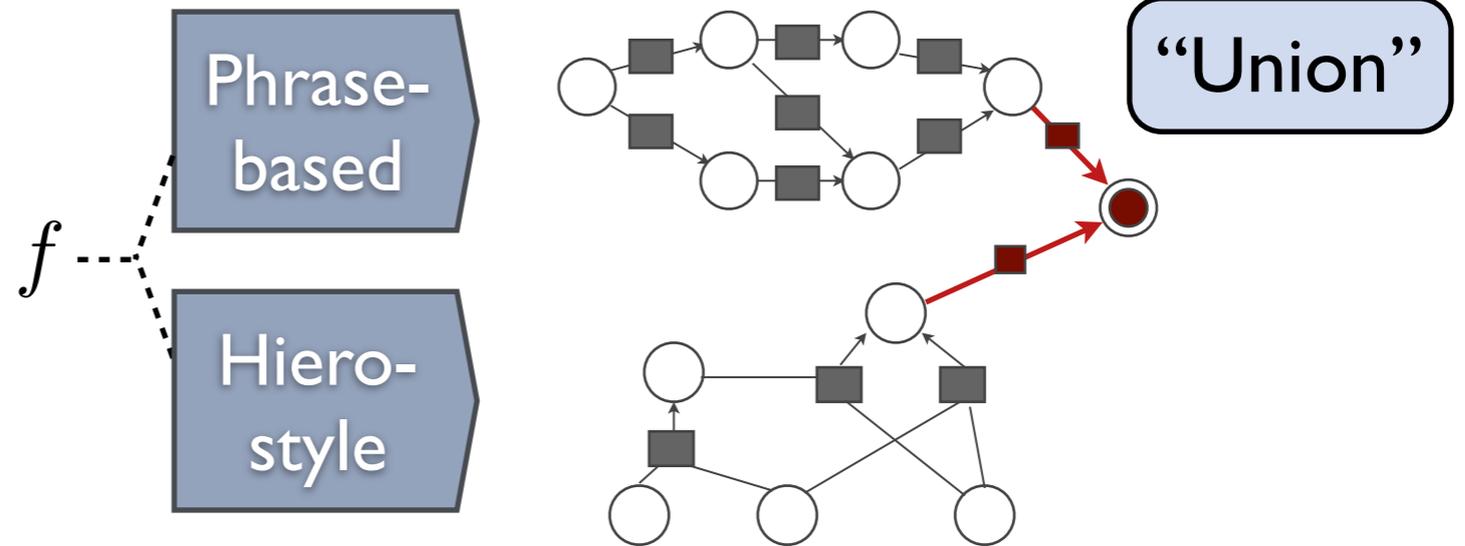
- Best max-derivation system
- Best single-system with consensus decoding
- Union + consensus decoding
- Union + consensus + model choice features
- Union with model-specific n-gram statistics



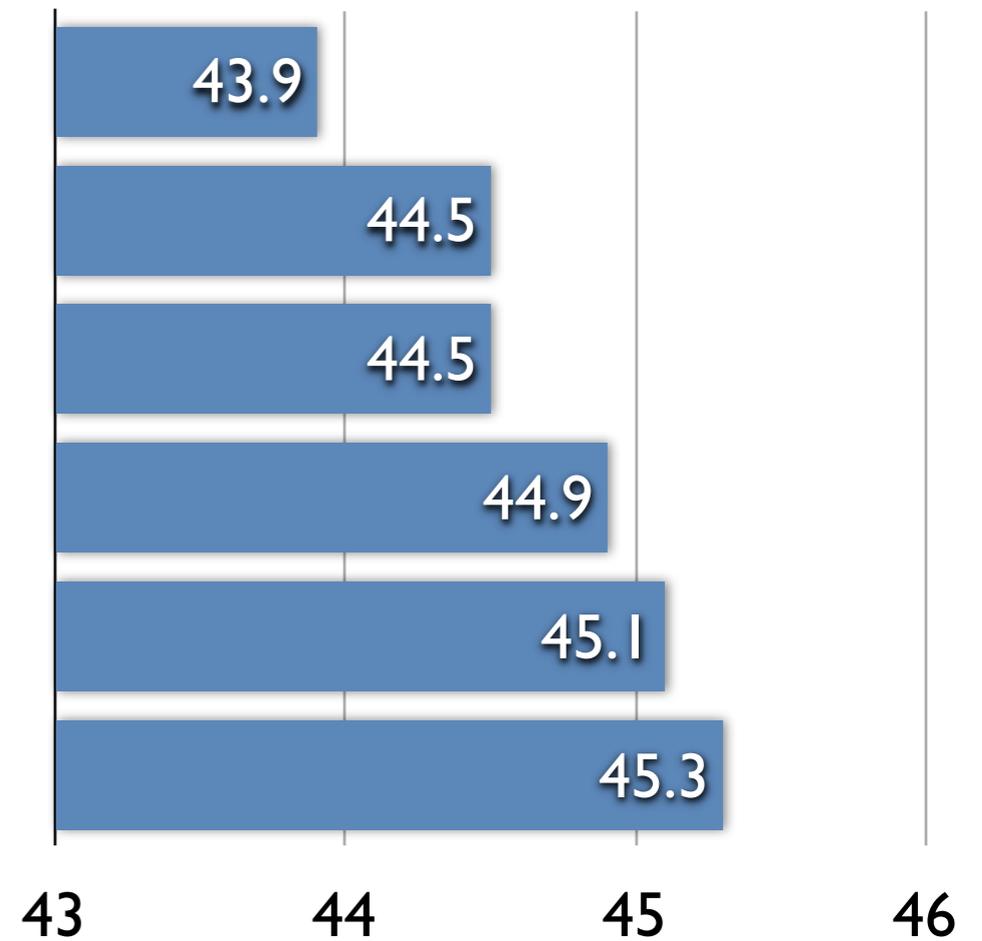
Sources of Improvement: Arabic-to-English

Google

- ▶ Do we need model choice features?
- ▶ Should n-gram statistics be model specific?



- Best max-derivation system
- Best single-system with consensus decoding
- Union + consensus decoding
- Union + consensus + model choice features
- Union with model-specific n-gram statistics
- Model-specific & union n-gram statistics



Outline



Consensus decoding review

Our model combination technique

Comparison to system combination

Outline

Consensus decoding review

Our model combination technique

Comparison to system combination



**The Final
Showdown**

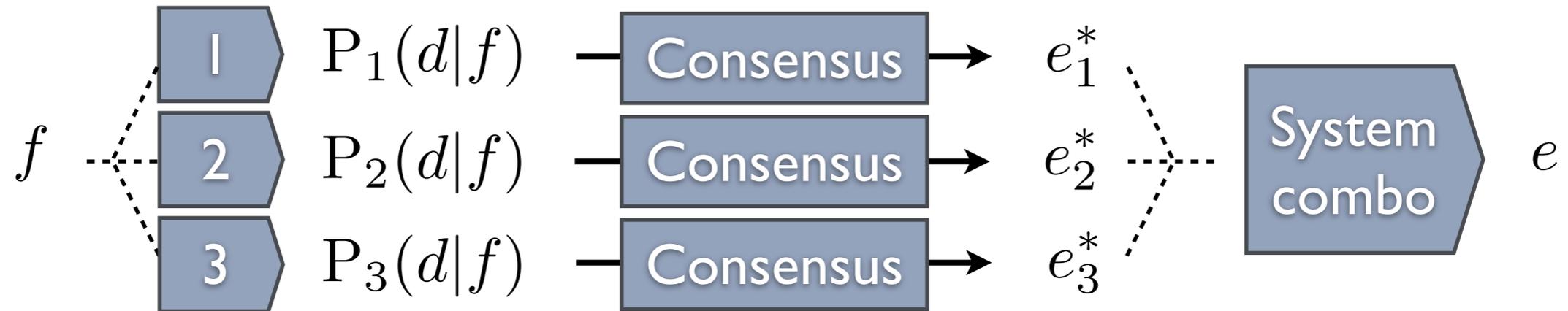
System Combination Baselines



Two established system combination methods [Macherey & Och, '07]

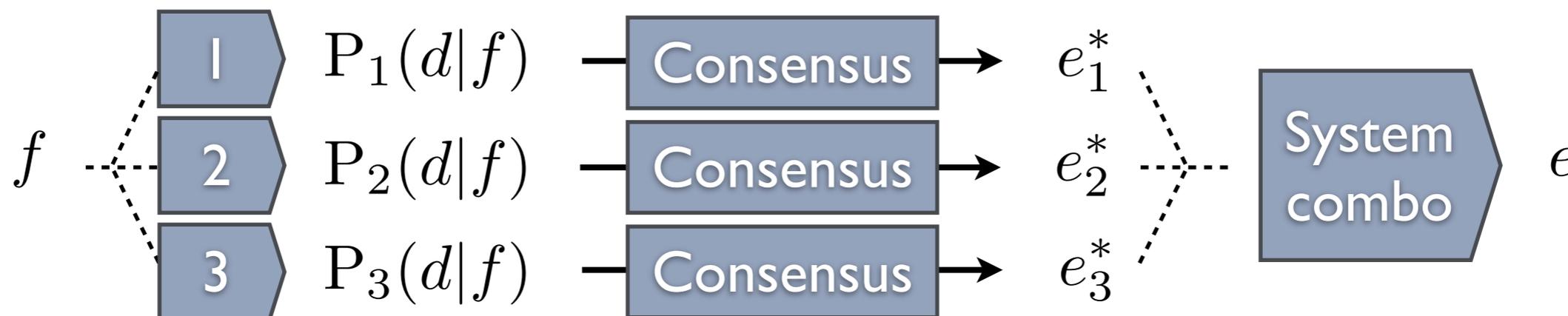
System Combination Baselines

Two established system combination methods [Macherey & Och, '07]



System Combination Baselines

Two established system combination methods [Macherey & Och, '07]



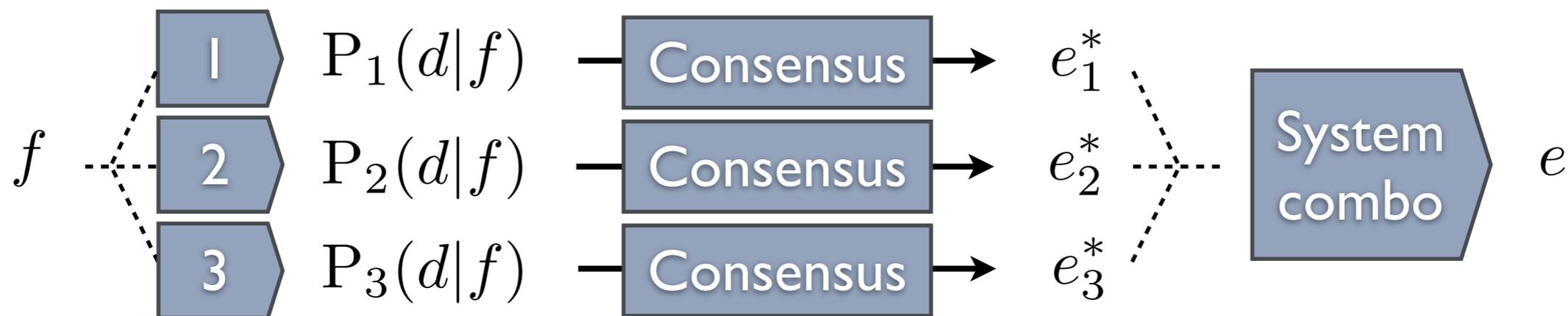
Sentence-level

Choose among system outputs using an MBR objective

$$e = \arg \max_{e \in \{e_1^*, \dots, e_k^*\}} \mathbb{E} [\text{BLEU}(e)]$$

System Combination Baselines

Two established system combination methods [Macherey & Och, '07]



Sentence-level

Choose among system outputs using an MBR objective

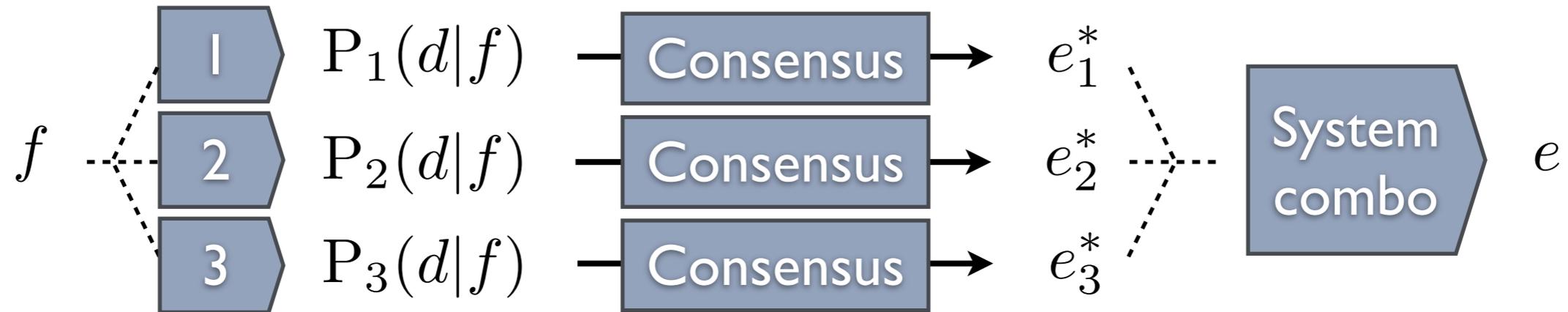
$$e = \arg \max_{e \in \{e_1^*, \dots, e_k^*\}} \mathbb{E} [\text{BLEU}(e)]$$

Word-level

Confusion network approach

System Combination Baselines

Two established system combination methods [Macherey & Och, '07]



Sentence-level

Choose among system outputs using an MBR objective

$$e = \arg \max_{e \in \{e_1^*, \dots, e_k^*\}} \mathbb{E} [\text{BLEU}(e)]$$

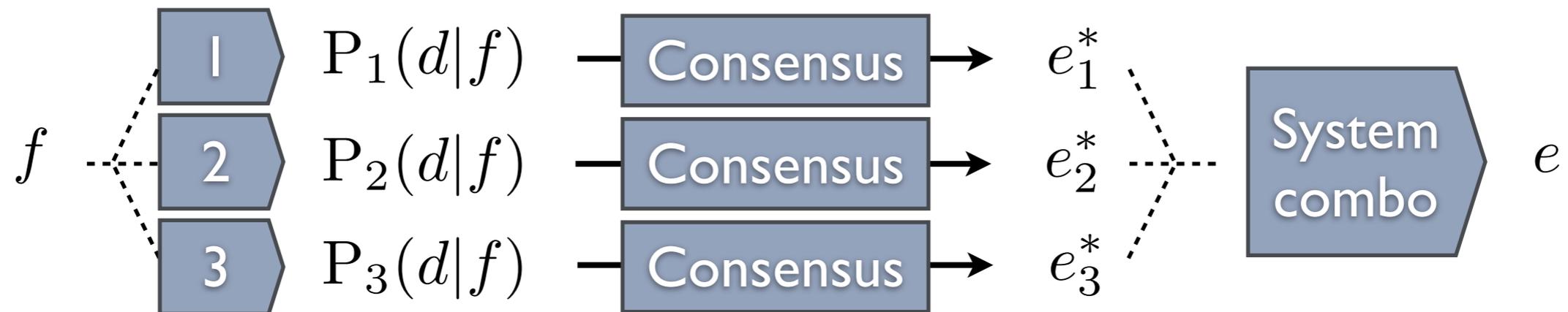
Word-level

Confusion network approach

- ▶ All e_i^* are aligned to a backbone $e_b \in \{e_1^*, \dots, e_k^*\}$

System Combination Baselines

Two established system combination methods [Macherey & Och, '07]



Sentence-level

Choose among system outputs using an MBR objective

$$e = \arg \max_{e \in \{e_1^*, \dots, e_k^*\}} \mathbb{E} [\text{BLEU}(e)]$$

Word-level

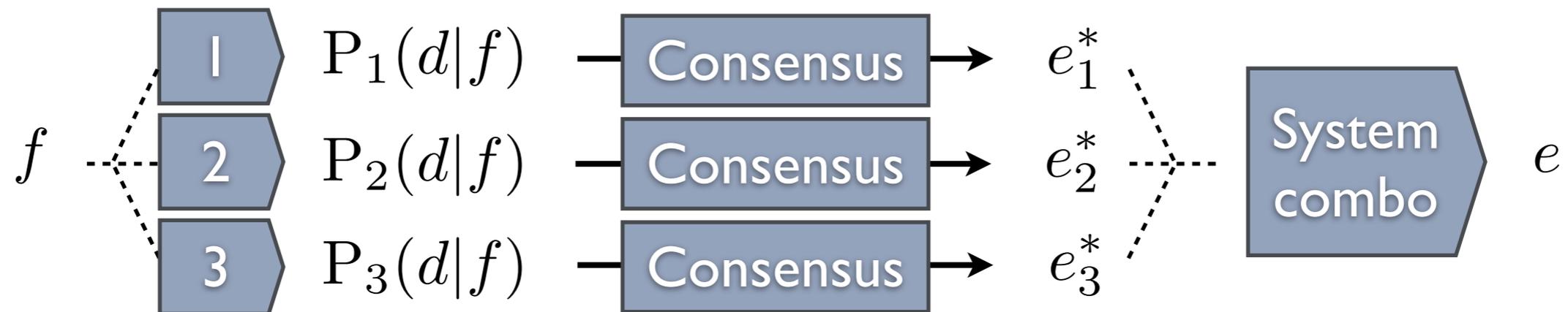
Confusion network approach

- ▶ All e_i^* are aligned to a backbone $e_b \in \{e_1^*, \dots, e_k^*\}$
- ▶ Sentences + alignments form a confusion network

System Combination Baselines

Google

Two established system combination methods [Macherey & Och, '07]



Sentence-level

Choose among system outputs using an MBR objective

$$e = \arg \max_{e \in \{e_1^*, \dots, e_k^*\}} \mathbb{E} [\text{BLEU}(e)]$$

Word-level

Confusion network approach

- ▶ All e_i^* are aligned to a backbone $e_b \in \{e_1^*, \dots, e_k^*\}$
- ▶ Sentences + alignments form a confusion network
- ▶ Output maximizes a consensus objective

Model versus System Combination

Google

Qualitative

Quantitative

Model versus System Combination

Google

Qualitative

- ▶ Single-system n-gram statistics are required in both methods

Quantitative

Model versus System Combination

Google

Qualitative

- ▶ Single-system n-gram statistics are required in both methods
- ▶ Model combination searches only once under a consensus objective

Quantitative

Model versus System Combination

Google

Qualitative

- ▶ Single-system n-gram statistics are required in both methods
- ▶ Model combination searches only once under a consensus objective
- ▶ Inter-hypothesis alignment problem exists only in system combination

Quantitative

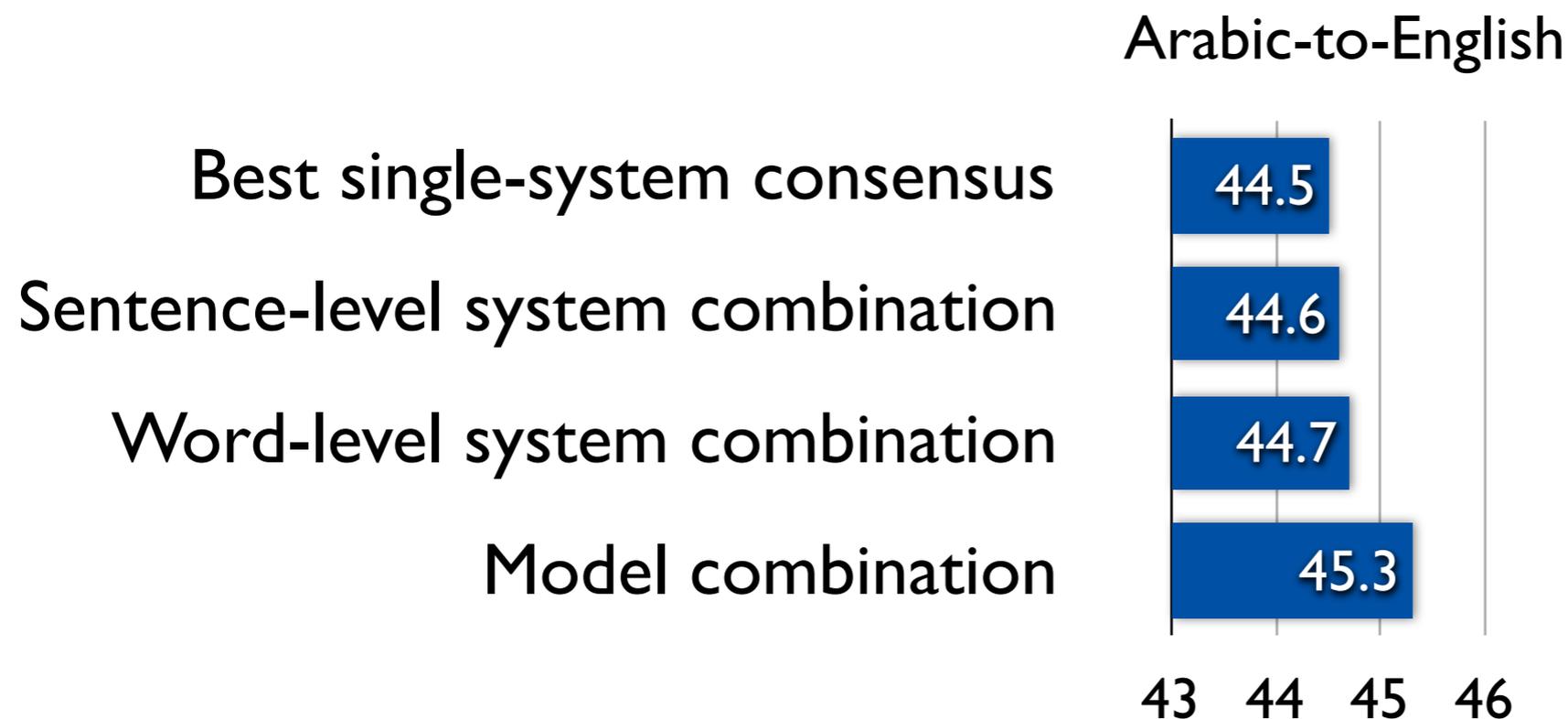
Model versus System Combination

Google

Qualitative

- ▶ Single-system n-gram statistics are required in both methods
- ▶ Model combination searches only once under a consensus objective
- ▶ Inter-hypothesis alignment problem exists only in system combination

Quantitative



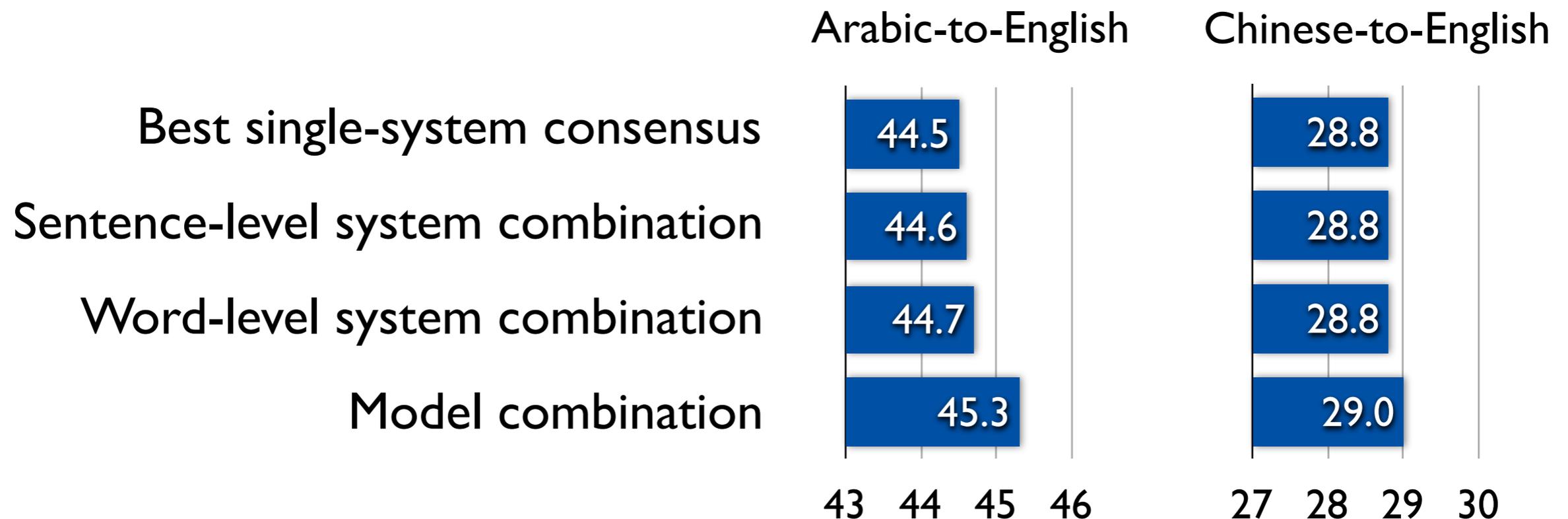
Model versus System Combination



Qualitative

- ▶ Single-system n-gram statistics are required in both methods
- ▶ Model combination searches only once under a consensus objective
- ▶ Inter-hypothesis alignment problem exists only in system combination

Quantitative



Conclusion

Google

Conclusion

Google

- ▶ Consensus decoding with a learned objective extends naturally to multiple models

Conclusion

Google

- ▶ Consensus decoding with a learned objective extends naturally to multiple models
- ▶ Model combination provides consensus *and* combination effects in a unified, distribution-driven objective

Conclusion

Google

- ▶ Consensus decoding with a learned objective extends naturally to multiple models
- ▶ Model combination provides consensus *and* combination effects in a unified, distribution-driven objective
- ▶ It outperforms a pipeline of consensus decoding followed by system combination, using less total computation

Conclusion

Google

- ▶ Consensus decoding with a learned objective extends naturally to multiple models
- ▶ Model combination provides consensus *and* combination effects in a unified, distribution-driven objective
- ▶ It outperforms a pipeline of consensus decoding followed by system combination, using less total computation

It's easy, it's clean, and it works

Conclusion



- ▶ Consensus decoding with a learned objective extends naturally to multiple models
- ▶ Model combination provides consensus *and* combination effects in a unified, distribution-driven objective
- ▶ It outperforms a pipeline of consensus decoding followed by system combination, using less total computation

It's easy, it's clean, and it works



Thanks!

A yellow starburst graphic with a black outline, containing the word "Thanks!" in a bold, black, sans-serif font. The starburst is positioned at the bottom right of the slide, overlapping the grey bar.