

# Efficient Minimum Error Rate Training and Minimum Bayes-Risk Decoding for Translation Hypergraphs and Lattices

Shankar Kumar, Wolfgang Macherey, Chris Dyer, and  
Franz Och

Presented By: Jeff Flanigan

April 18, 2012

# Overview

- Hypergraph MERT
- Efficient Minimum Bayes-Risk (MBR) Decoding for Lattices
- MBR Decoding on Hypergraphs
- Combine all three
- Results

# MERT Review

Decoder rule: 
$$\hat{E}(F_s; \lambda_1^M) = \arg \max_E \left\{ \sum_{m=1}^M \lambda_m h_m(E, F_s) \right\}$$

MERT Tuning: 
$$\hat{\lambda}_1^M = \arg \min_{\lambda_1^M} \left\{ \sum_{s=1}^S \text{Err}(R_s, \hat{E}(F_s; \lambda_1^M)) \right\}$$

$F_s$ : source

$R_s$ : reference translation

$h_m$ : feature functions

$\lambda_m$ : feature weights

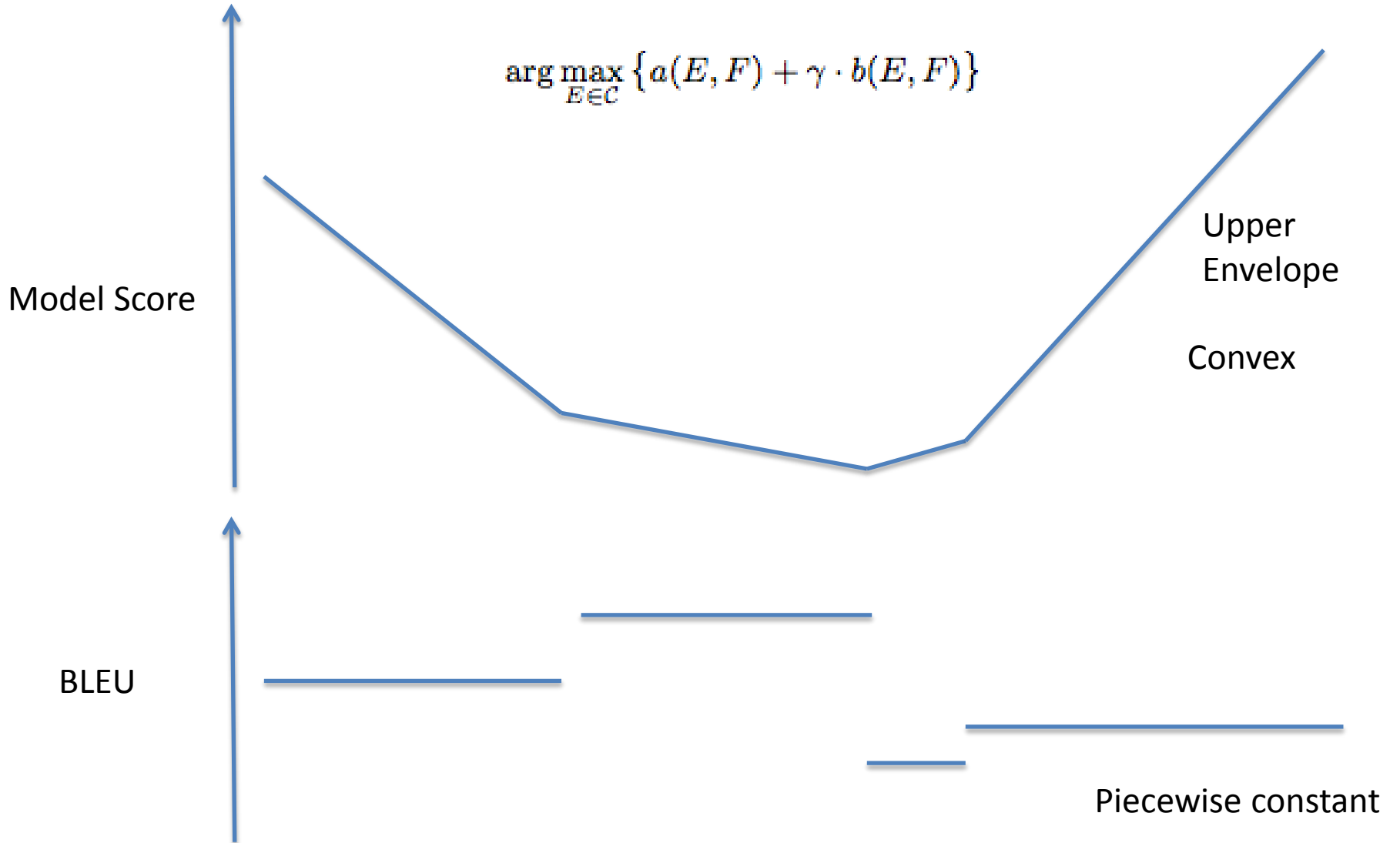
# MERT Review

Decoder rule:  $\hat{E}(F_s; \lambda_1^M) = \arg \max_E \left\{ \sum_{m=1}^M \lambda_m h_m(E, F_s) \right\}$

Choose direction  $d_1^M$   $\gamma = \lambda_1^M + \gamma \cdot d_1^M$

$$\begin{aligned} \hat{E}(F; \gamma) &= \arg \max_{E \in \mathcal{C}} \left\{ (\lambda_1^M + \gamma \cdot d_1^M)^\top \cdot h_1^M(E, F) \right\} \\ &= \arg \max_{E \in \mathcal{C}} \left\{ \underbrace{\sum_m \lambda_m h_m(E, F)}_{=a(E, F)} + \gamma \cdot \underbrace{\sum_m d_m h_m(E, F)}_{=b(E, F)} \right\} \\ &= \arg \max_{E \in \mathcal{C}} \{ a(E, F) + \gamma \cdot b(E, F) \} \end{aligned}$$

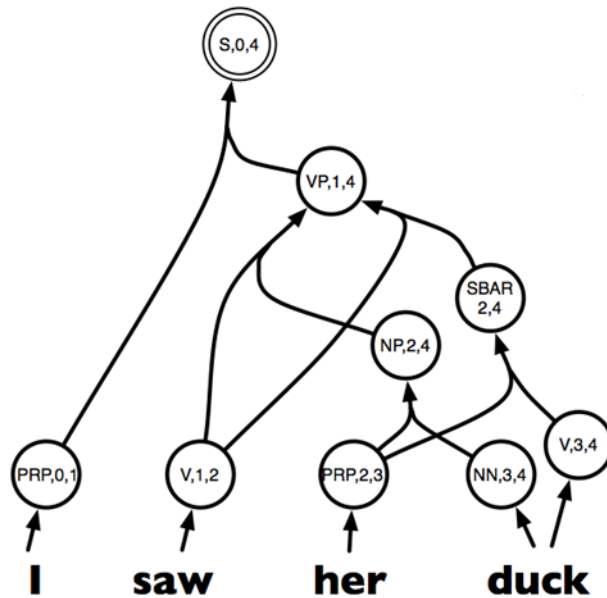
# MERT Review



# Hypergraphs

Hypergraph has nodes and edges  $H = \langle V, E \rangle$

Edges directed: go from multiple tails to single head

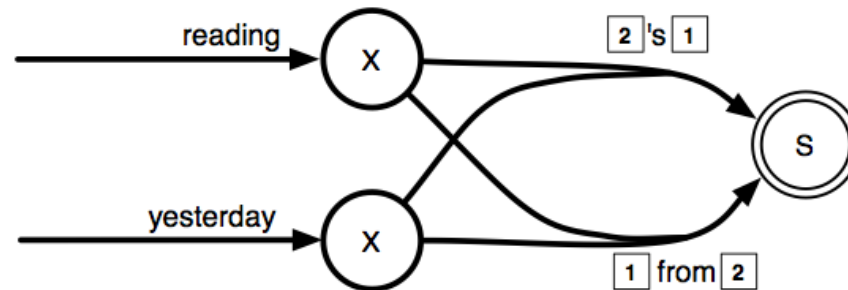


# Translation Hypergraphs

Each edge labeled with a rule

Nodes = nonterminals

Path = translation hypothesis



$\left\{ \left( \text{yesterday 's reading} \right), \right.$   
 $\left. \left( \text{reading from yesterday} \right) \right\}$

# Hypergraph MERT

Repeat until convergence:

- Pick a direction

- Efficiently calculate upper envelope over entire lattice

- Line search for best BLEU score (for entire devset)



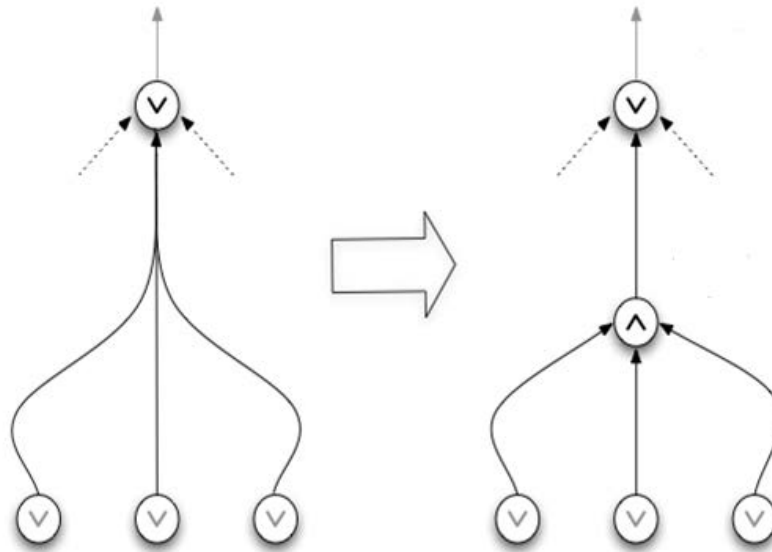
# Hypergraph MERT

## Computing the upper envelope

Step 1: Convert hypergraph to regular graph

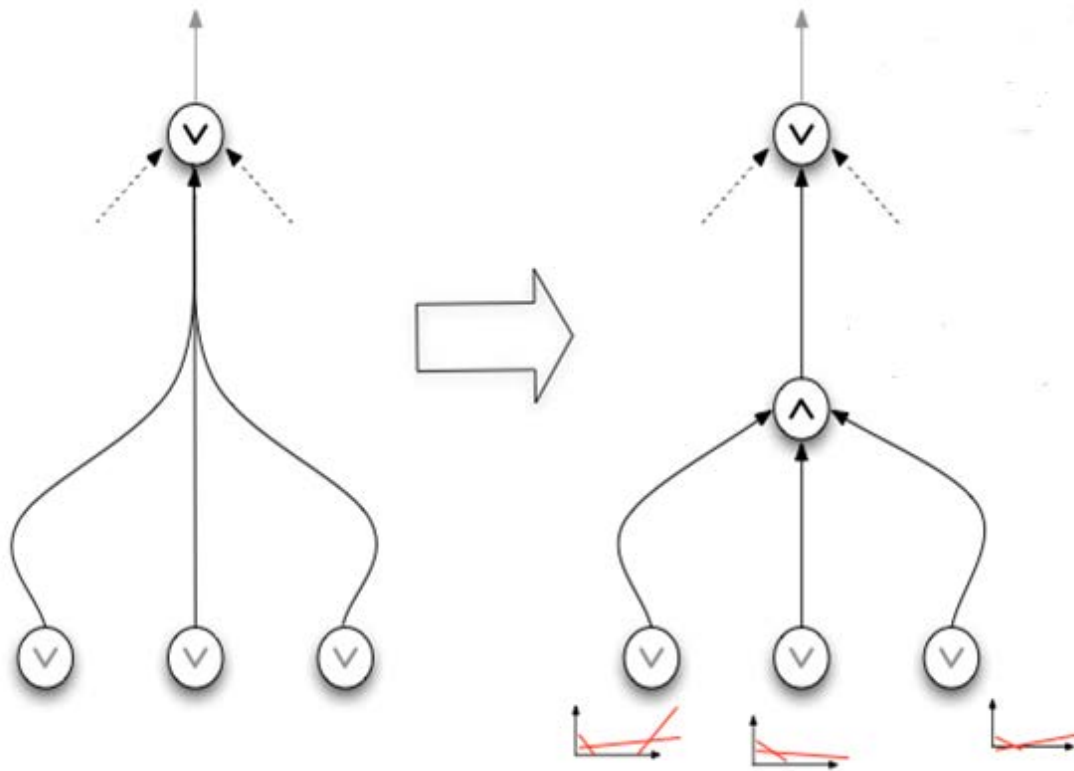
Nodes  $\rightarrow$   $\vee$  type nodes

Edges  $\rightarrow$   $\wedge$  type nodes



# Hypergraph MERT

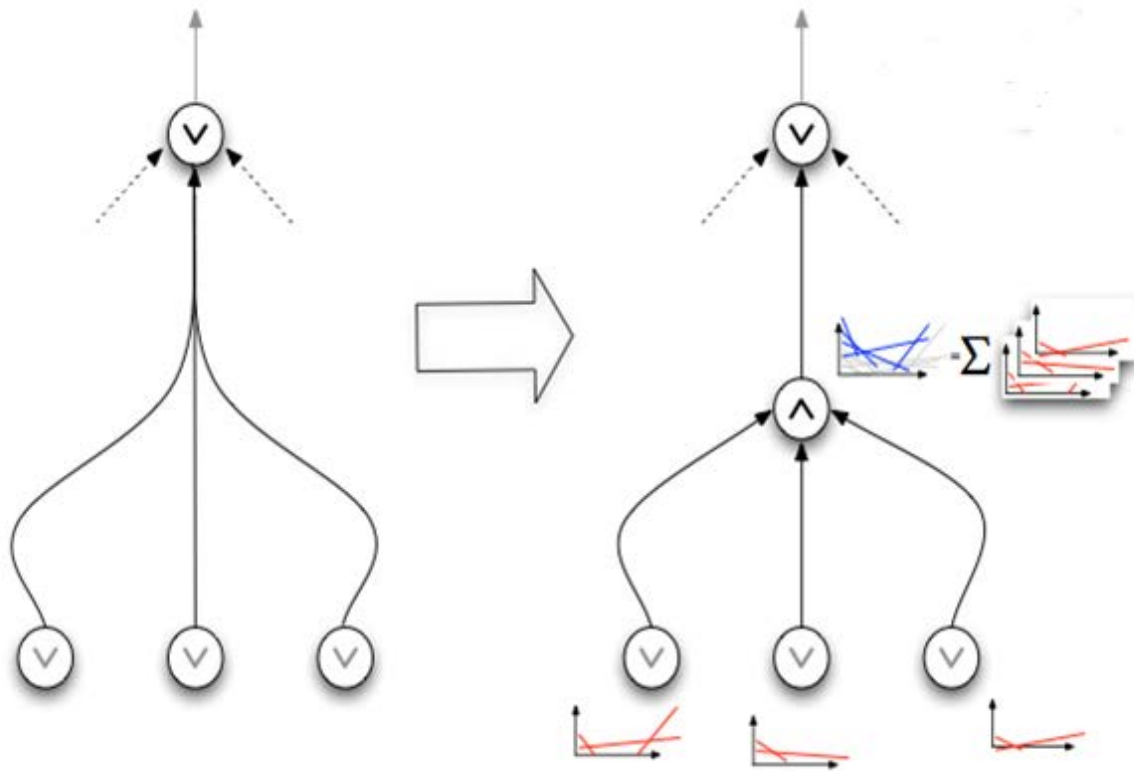
Step 2: Propagate upper envelope up to the root



# Hypergraph MERT

Step 2: Propagate upper envelope up to the root

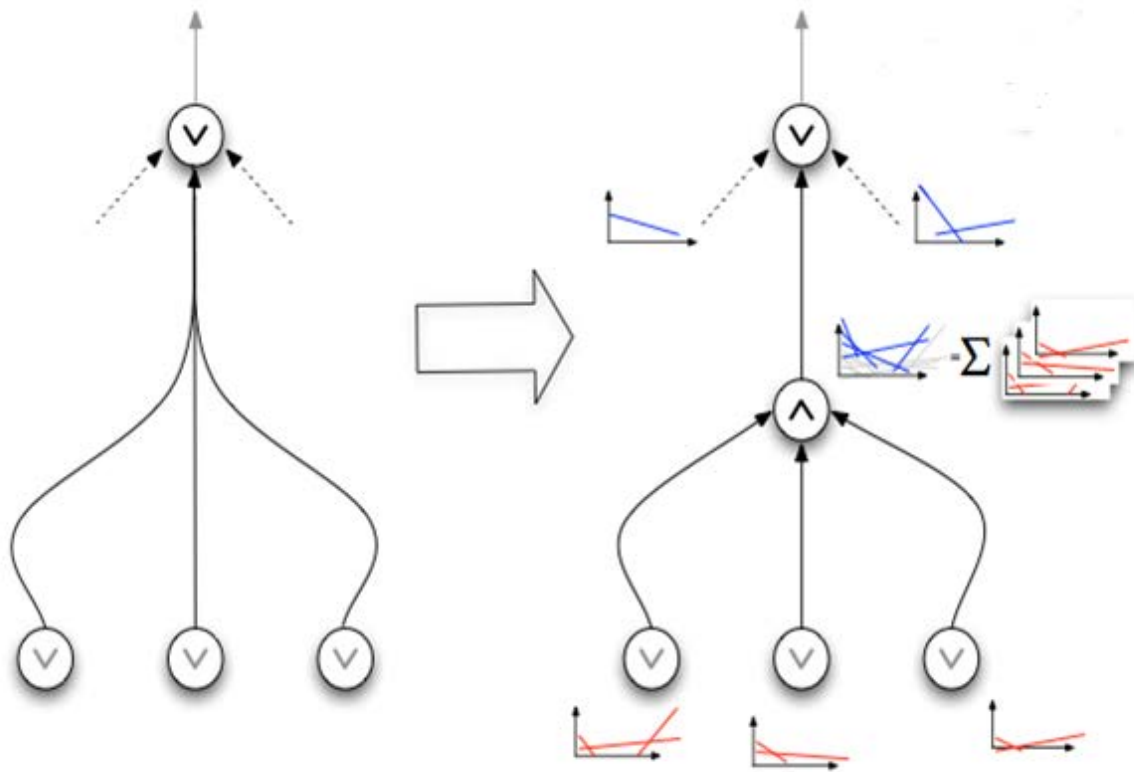
$\wedge$  nodes: Sum (also include rule score)



# Hypergraph MERT

Step 2: Propagate upper envelope up to the root

$\wedge$  nodes: Sum (also include rule score)

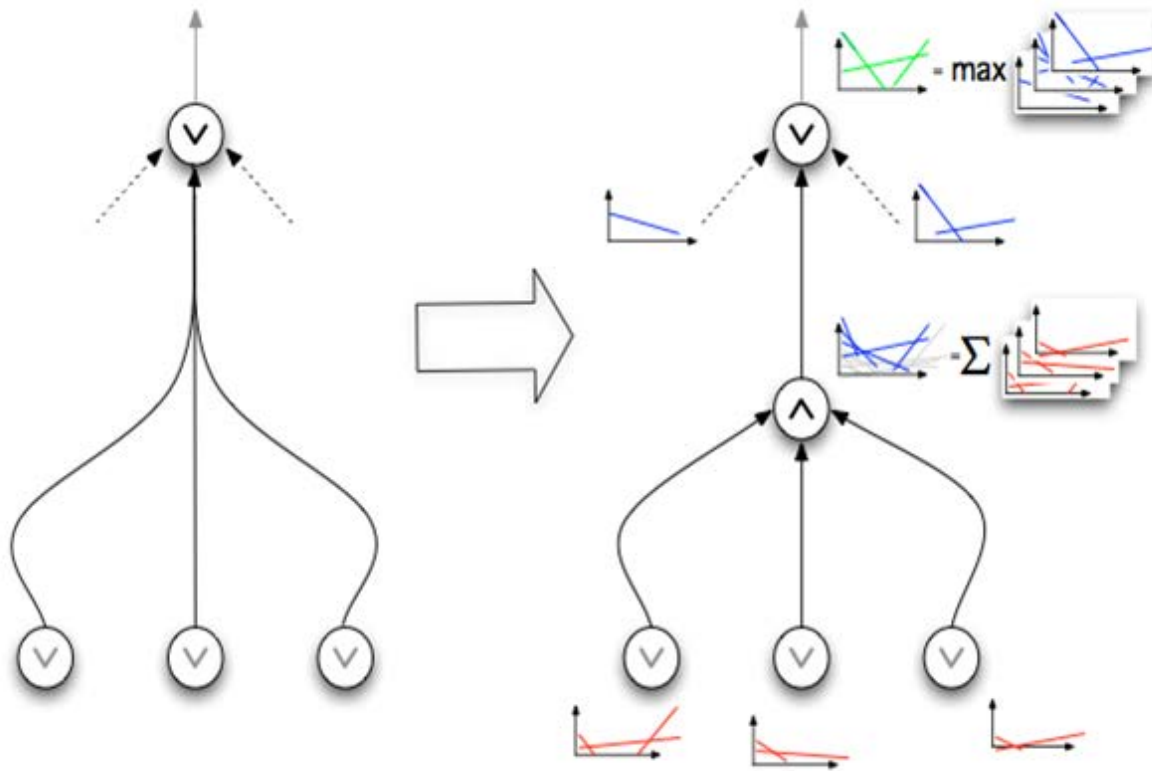


# Hypergraph MERT

Step 2: Propagate upper envelope up to the root

$\wedge$  nodes: Sum (also include rule score)

$\vee$  nodes: Max



# Minimum Bayes-Risk Decoding

Decoding Rule:

$$\hat{E} = \operatorname{argmin}_{E' \in \mathcal{G}} \sum_{E \in \mathcal{G}} L(E, E') P(E|F)$$

Minimize expected loss under  
probability model  $P(E|F)$

# Minimum Bayes-Risk Decoding

Decoding Rule:

$$\hat{E} = \operatorname{argmin}_{E' \in \mathcal{G}} \sum_{E \in \mathcal{G}} L(E, E') P(E|F)$$

1 if w is in E  
0 otherwise

-Loss = log(BLEU)  $\approx$

$$G(E, E') = \theta_0 |E'| + \sum_w \theta_{|w|} \#_w(E') \delta_w(E)$$

# of times w is  
in E'

# Minimum Bayes-Risk Decoding

Decoding Rule:

$$\hat{E} = \operatorname{argmin}_{E' \in \mathcal{G}} \sum_{E \in \mathcal{G}} L(E, E') P(E|F)$$

1 if w is in E  
0 otherwise

-Loss = log(BLEU)  $\approx$

$$G(E, E') = \theta_0 |E'| + \sum_w \theta_{|w|} \#_w(E') \delta_w(E)$$

Sum over E  
 $P(E|F) = 1$

# of times w is  
in E'

Plug in  $\Rightarrow$

$$\hat{E} = \operatorname{argmax}_{E' \in \mathcal{G}} \theta_0 |E'| + \sum_w \theta_{|w|} \#_w(E') p(w|\mathcal{G})$$

where

$$p(w|\mathcal{G}) = \sum_{E \in \mathcal{G}} 1_w(E) P(E|F)$$

“Posterior probability of  
n-gram w in the lattice G”



# Efficient MBR on Lattices

Rewrite

$$p(w|\mathcal{G}) = \sum_{E \in \mathcal{G}} 1_w(E) P(E|F)$$

$$p(w|\mathcal{G}) = \sum_{E \in \mathcal{G}} \sum_{e \in E} f(e, w, E) P(E|F)$$

where

$$f(e, w, E) = \begin{cases} 1 & w \in e, p(e|\mathcal{G}) > p(e'|\mathcal{G}), \\ & e' \text{ precedes } e \text{ on } E \\ 0 & \text{otherwise} \end{cases}$$

Count w  
only once in  
hypothesis

Approximate  $f$  with  $f^*$  = indicates edge containing  $w$  that has highest arc probability  $p(e|\mathcal{G})$   
 **$f^*$  can be calculated independent of path  $\Leftarrow$  efficient**

$$p(w|\mathcal{G}) = \sum_{E \in \mathcal{G}} \sum_{e \in E} f^*(e, w, \mathcal{G}) P(E|F) \quad (7)$$

$$= \sum_{e \in \mathcal{E}} 1_{w \in e} f^*(e, w, \mathcal{G}) \sum_{E \in \mathcal{G}} 1_E(e) P(E|F)$$

$e$  ranges over all edges,  
pull out of sum

$$= \sum_{e \in \mathcal{E}} 1_{w \in e} f^*(e, w, \mathcal{G}) P(e|\mathcal{G}),$$

$P(e|\mathcal{G})$  is posterior prob of lattice edge

# Efficient MBR on Lattices

---

## Algorithm 3 MBR Decoding on Lattices

---

- 1: Sort the lattice nodes topologically.
  - 2: Compute backward probabilities of each node.
  - 3: Compute posterior prob. of each  $n$ -gram:
  - 4: **for** each edge  $e$  **do**
  - 5:     Compute edge posterior probability  $P(e|\mathcal{G})$ .
  - 6:     Compute  $n$ -gram posterior probs.  $P(w|\mathcal{G})$ :
  - 7:     **for** each  $n$ -gram  $w$  introduced by  $e$  **do**
  - 8:         Propagate  $n - 1$  gram suffix to  $h_e$ .
  - 9:         **if**  $p(e|\mathcal{G}) > \text{Score}(w, T(e))$  **then**
  - 10:             Update posterior probs. and scores:  
                $p(w|\mathcal{G}) += p(e|\mathcal{G}) - \text{Score}(w, T(e))$ .  
                $\text{Score}(w, h_e) = p(e|\mathcal{G})$ .
  - 11:             **else**
  - 12:                  $\text{Score}(w, h_e) = \text{Score}(w, T(e))$ .
  - 13:             **end if**
  - 14:     **end for**
  - 15: **end for**
  - 16: Assign scores to edges (given by Equation 3).
  - 17: Find best path in the lattice (Equation 3).
- 

Score( $w, t$ ) is  
highest probability  
of paths that  
terminate on  $t$  and  
contain  $n$ -gram  $w$

# Efficient MBR on Hypergraphs

---

## Algorithm 4 MBR Decoding on Hypergraphs

---

- 1: Sort the hypergraph nodes topologically.
  - 2: Compute inside probabilities of each node.
  - 3: Compute posterior prob. of each hyperedge  $P(e|\mathcal{G})$ .
  - 4: Compute posterior prob. of each  $n$ -gram:
  - 5: **for** each hyperedge  $e$  **do**
  - 6:     Merge the  $n$ -grams on the tail nodes  $T(e)$ . If the same  $n$ -gram is present on multiple tail nodes, keep the highest score.
  - 7:     Apply the rule on  $e$  to the  $n$ -grams on  $T(e)$ .
  - 8:     Propagate  $n - 1$  gram prefixes/suffixes to  $h_e$ .
  - 9:     **for** each  $n$ -gram  $w$  introduced by this hyperedge **do**
  - 10:         **if**  $p(e|\mathcal{G}) > \text{Score}(w, T(e))$  **then**
  - 11:              $p(w|\mathcal{G}) += p(e|\mathcal{G}) - \text{Score}(w, T(e))$
  - 12:              $\text{Score}(w, h_e) = p(e|\mathcal{G})$
  - 13:             **else**
  - 14:                  $\text{Score}(w, h_e) = \text{Score}(w, T(e))$
  - 15:             **end if**
  - 16:     **end for**
  - 17:     Assign scores to hyperedges (Equation 3).
  - 18: **end for**
  - 19: Find best path in the hypergraph (Equation 3).
- 

Essentially the same as for lattices. Need to propagate  $n$ -gram prefixes and suffices

# MERT for MBR Parameter Tuning

Use MERT to tune  $\theta_i$ 's

$$\hat{E} = \operatorname{argmax}_{E' \in \mathcal{G}} \theta_0 |E'| + \sum_w \theta_{|w|} \#_w(E') p(w|\mathcal{G})$$

Also include additional feature  $g_{N+1}(E,F) = \text{original decoder cost (MAP translation)}$

# Evaluation

NIST training data

Dataset	# of sentences	
	aren	zhen
dev	1797	1664
nist02	1043	878
nist03	663	919

	Avg. Runtime/sent [msec]			
	(Macherey 2008)		Suggested Alg.	
	aren	zhen	aren	zhen
phrase lattice	8.57	7.91	10.30	8.65
hypergraph	-	-	8.19	8.11

Table 2: Average time for computing envelopes.

Pruned hypergraph comparable running time

	BLEU (%)				Avg. time (ms.)
	aren		zhen		
	nist03	nist02	nist03	nist02	
MAP	54.2	64.2	40.1	39.0	-
<i>N</i> -best MBR	54.3	64.5	40.2	39.2	3.7
Lattice MBR					
FSAMBR	54.9	65.2	40.6	39.5	3.7
LatMBR	54.8	65.2	40.7	39.4	0.2

Table 3: Lattice MBR for a phrase-based system.

Lattice MBR 20x faster than FSAMBR

	BLEU (%)				Avg. time (ms.)
	aren		zhen		
	nist03	nist02	nist03	nist02	
Hiero					
MAP	52.8	62.9	41.0	39.8	-
<i>N</i> -best MBR	53.2	63.0	41.0	40.1	3.7
HGMBR	53.3	63.1	41.0	40.2	0.5
SAMT					
MAP	53.4	63.9	41.3	40.3	-
<i>N</i> -best MBR	53.8	64.3	41.7	41.1	3.7
HGMBR	54.0	64.6	41.8	41.1	0.5

Table 4: Hypergraph MBR for Hiero/SAMT systems.

Hypergraph MBR 7x faster than *N*-best MBR

# Evaluation

System	BLEU (%)			
	MAP	MBR		
		default	mert-b	mert+b
aren.pb	54.2	54.8	54.8	54.9
aren.hier	52.8	53.3	53.5	53.7
aren.samt	53.4	54.0	54.4	54.0
zhen.pb	40.1	40.7	40.7	40.9
zhen.hier	41.0	41.0	41.0	41.0
zhen.samt	41.3	41.8	41.6	41.7

Not much gain over default parameter settings

Table 5: MBR Parameter Tuning on NIST systems

MBR wrt. MAP	default	mert-b	mert+b
# of gains	18	22	26
# of no-changes	9	5	8
# of drops	12	12	5

Gain over default parameter settings

Table 6: MBR on Multi-language systems.