OPTIMAL SEARCH FOR MINIMUM ERROR RATE TRAINING

Authors: Michel Galley & Chris Quirk (Microsoft Research)

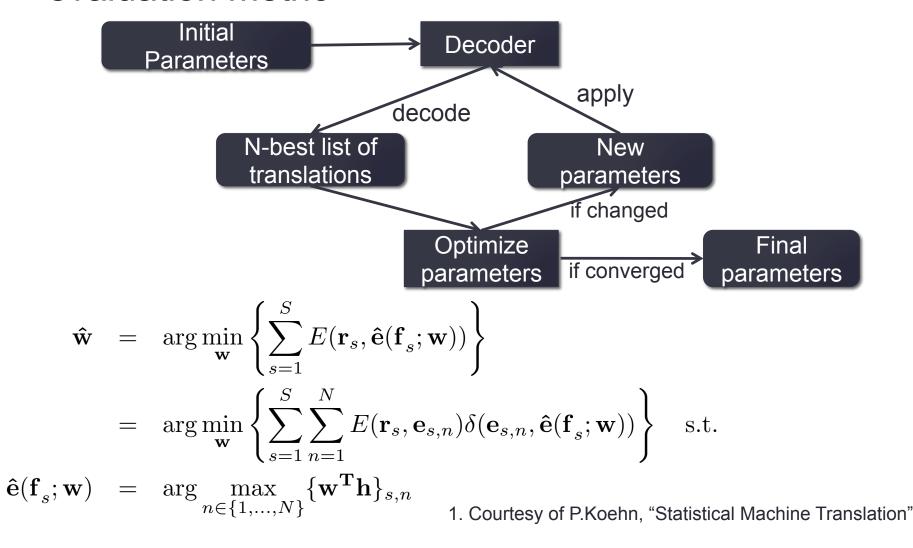
Presenter: Avneesh Saluja

The Claim

- Och's MERT is not exact
 - I will provide a brief MERT review
- This paper: exact search of MERT search space via linear programming
 - Concurrent optimization of all dimensions
- How does this perform vs. Och's MERT?

MERT for MT: Primer¹

Tune parameter weights to directly optimize evaluation metric



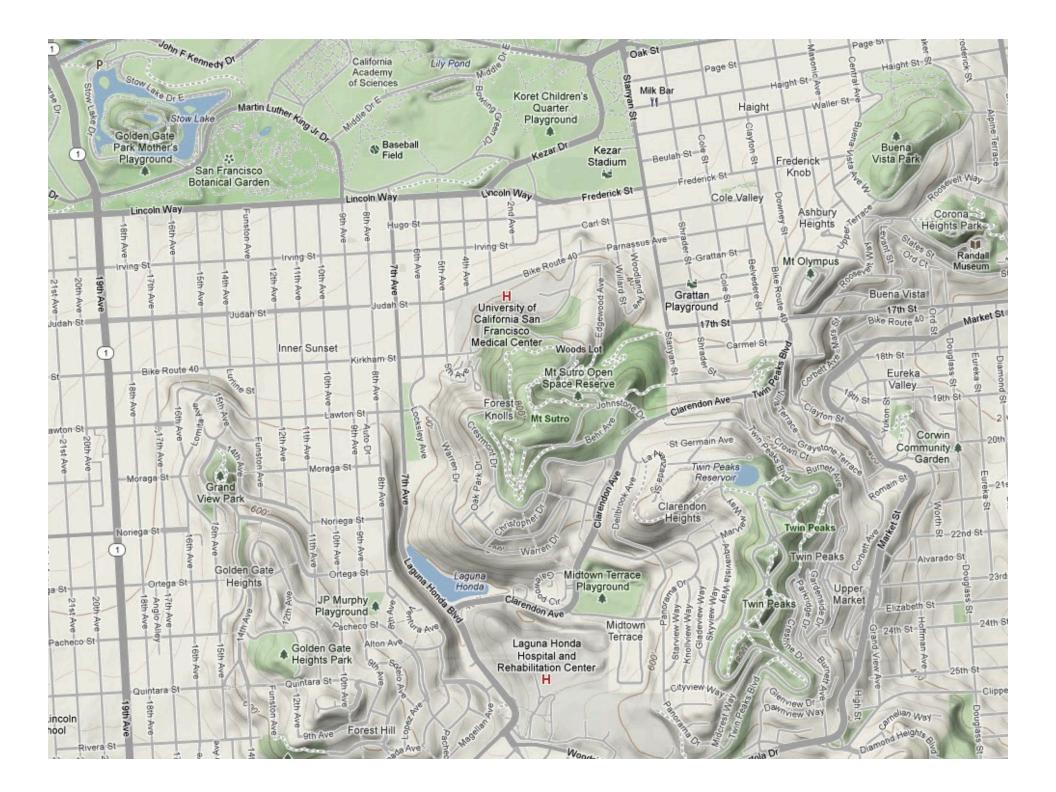
Naïve MERT

- Non-convex, unsmooth: Powell search
 - Multiple starting points used
- Log-linear formulation: $p(x) = \exp\left\{\sum_{i=1}^{N} \lambda_i h_i(x)\right\}$
- Best translation:

$$x^*(\lambda_1, \dots, \lambda_n) = \arg\max_{x} \exp\left\{\sum_{i=1}^{N} \lambda_i h_i(x)\right\} \Rightarrow$$

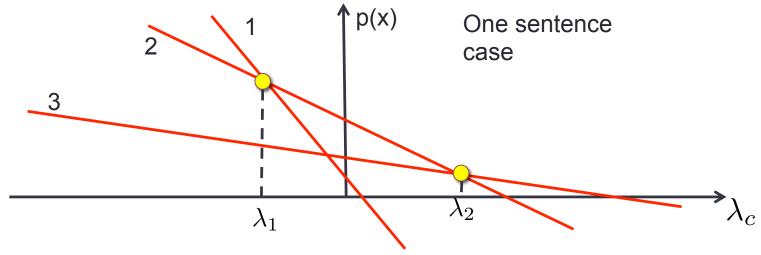
$$x^*(\lambda_c) = \arg\max_{x} \exp\left\{\lambda_c h_c(x) + u(x)\right\} \text{ where } u(x) = \sum_{i \neq c} \lambda_i h_i(x)$$

- How to explore parameter space?
 - Grid: trade-off between speed & accuracy
 - Finite approximation



Och's Trick

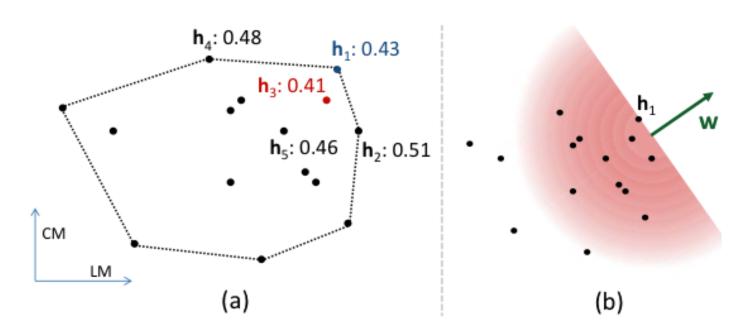
 Translation ranking only changes at intersections of translation lines!



- Multiple sentences:
 - aggregate intersections across sentences
 - For each intersection point, compute error (re-rank and select 1-best)
 - Return interval with best error score

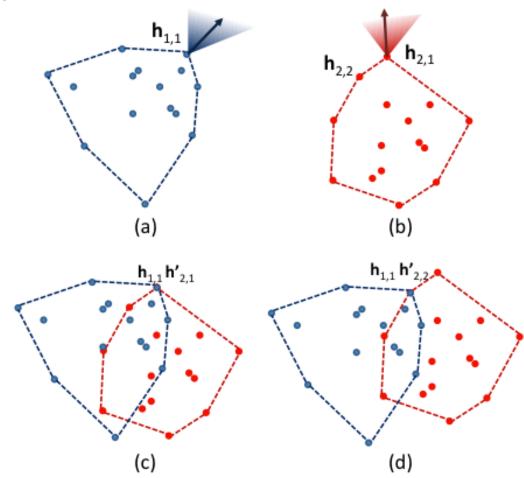
LP-MERT: one sentence case

- Och's MERT: great, but optimizing one parameter at a time?
 - Search over a larger subspace of parameter combos, not just line search
- Build convex hull of n-best list, iterate through extreme points
 - CH construction algorithms exponential in dimension
- Resort to LP with interior point methods (poly in dimension) to find extreme points



And more than 1 sentence?

- LP Formulation: return 0 if interior point
- Naïve approach:
 enumerate all possible
 hypothesis combinations
 across all sentences
- Smarter approach: merging convex hulls to maintain convexity
- Extreme point determination: O(NS) # points vs. $O(N^S)$



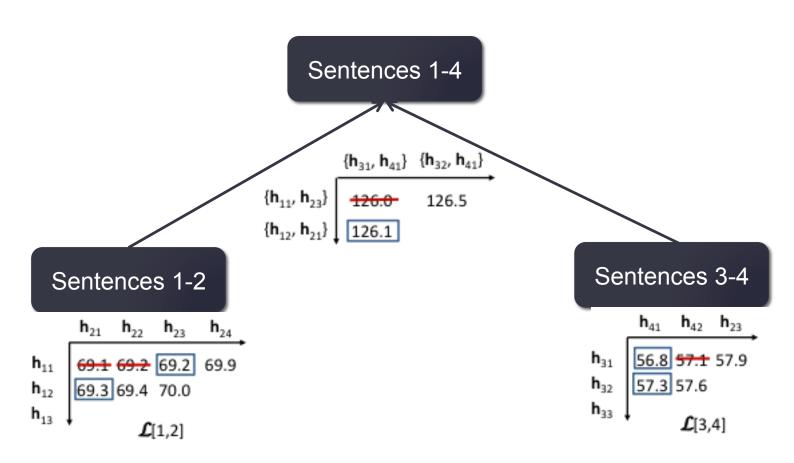
Other Tricks and Speedups

- Takes O(NS) points to determine if a point is extreme. Need to do this for $O(N^S)$ possible combinations
- Trick 1: lazy enumeration (ordering of combos)
- However, not quite enough

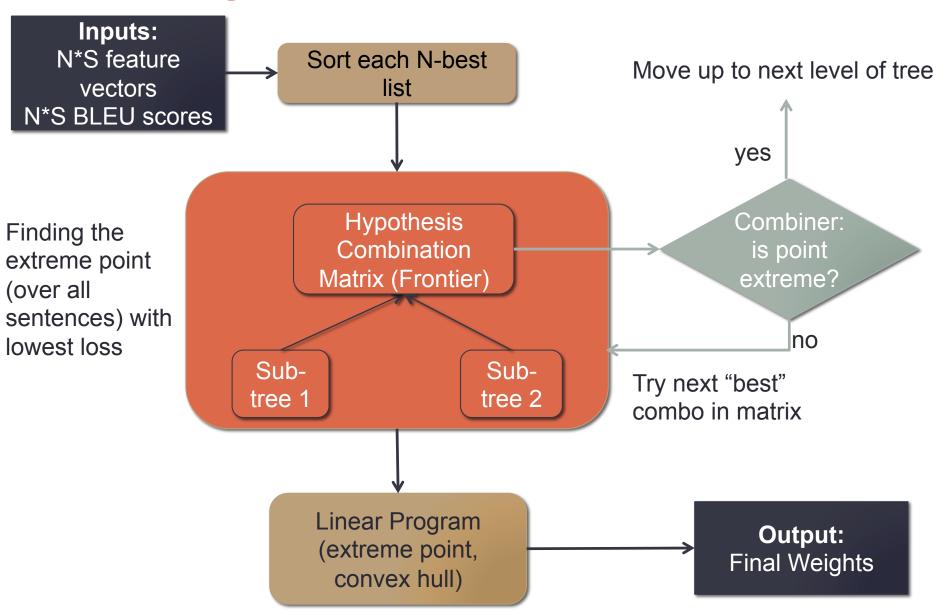
Sentence 1 Sentence 2 Sentence 3 Sentence 4 — Pizza with coke — That's sick bro! Where are Hello Pizza with drink —— That brother ill you? Good Day — Pizza coke That is ill Where is you? Hello brother Where are 1? Hi Pizza with coke Brother ill that Who are you?

Binary lazy enumeration

Use divide-and-conquer:

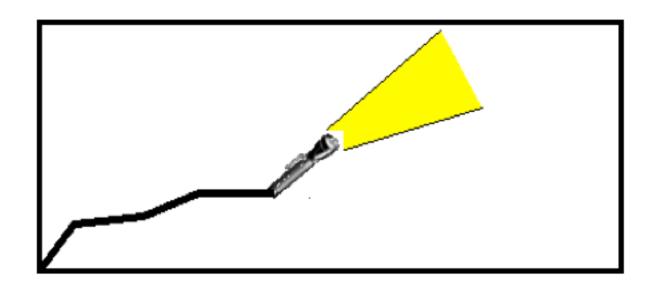


Final Algorithm



Approximations that we need

- Cosine similarity check (with reference vector): $cos(\hat{\mathbf{w}}, \mathbf{w_0}) \ge t$
- Beam search: prune with respect to current best parameter vector (when combining, check model score)



Experimental Setup

- Tree-to-string model
- 13 features in total
 - Standard PM and LM features, re-ordering, function word insertion/ deletion, insertion/deletion counts, target length
- N-best size = 100
 - Same combined N-best lists
- WMT 2010 English → German (1.6 million sentence pairs)
 - 2009 test: tuning
 - 2010 test: test
 - One reference translation

The D&C Speedup

length	tested comb.	total comb.	order
8	639,960	1.33×10^{20}	$O(N^8)$
4	134,454	2.31×10^{10}	$O(2N^4)$
2	49,969	430,336	$O(4N^{2})$
1	1,059	2,624	O(8N)

Table 1: Number of tested combinations for the experiments of Fig. 5. LP-MERT with S=8 checks only 600K full combinations on average, much less than the total number of combinations (which is more than 10^{20}).

Dependence on dimension

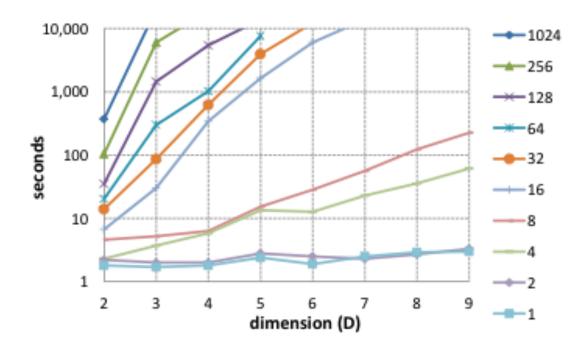


Figure 6: Effect of the number of features (runtime on 1 CPU of a modern computer). Each curve represents a different number of tuning sentences.

Cosine Similarity Approximation

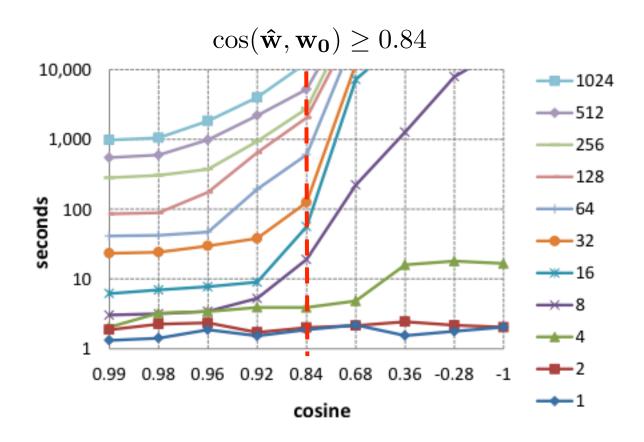
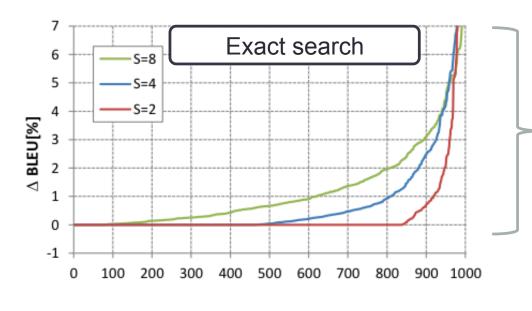


Figure 7: Effect of a constraint on w (runtime on 1 CPU).

Comparison with 1D-MERT

Assuming LP-MERT finds the "global optimum", 'for S=4, [Powell] makes search errors in 90% of the cases, despite using 20 random starting points'



With beam (size 1000)	32	64	128	256	512	1024		As S increas
1D-MERT	I							gap bet
our work							- 1	1D-ME
	+2.32	+1.59	+1.29	+0.98	+0.56	+0.23		and LP

Table 2: BLEUn4r1[%] scores for English-German on WMT09 for tuning sets ranging from 32 to 1024 sentences.

As S increases, the gap between 1D-MERT and LP-MERT decreases

As S

increases, the

gap between

and LP-MERT

1D-MERT

increases

Summary

- End-to-end evaluation (with beam approx. for LP-MERT)
 - Tuning: 0.24 BLEU difference
 - Test: 0.17 BLEU difference
- Exact multi-dimensional MERT
 - LP at the core
 - Divide-and-conquer, lazy enumeration
- Polynomial in dimension, N-best list size
- Exponential in number of sentences
- Approximations used to limit running time
- Bold approach to tackle difficult problem

Questions & concerns that I had...

- End-to-end results do not look significant
- Additional language pairs/datasets would be nice
- As S increases, does the over-performance diminish?
- What can we do to make this algorithm poly(S)?
- LP-MERT + hypergraph MERT → towards MERT
 2.0?
- Is direct cost optimization the way forward?

Thank you!