

# Dual, Monolingual, Cross-Entropy-Delta Filtering of Noisy Parallel Data



Amittai Axelrod  
DiDi Labs, Los Angeles  
amittai@didiglobal.com

Anish Kumar

Steve Sloto

- **One laptop, 3 hours.**
- **Nothing pretrained**
- **No parallel data**
- **No multilingual tools**
- Uses relative informativeness of data
- Planet-friendly, and surprisingly decent!

(...considering you have no parallel data)

Inspiration:

Filtering with Dual Conditional Cross-Entropy

[Junczys-Dowmunt, 2018]

Dual MT systems:

parallel inputs → identical likelihoods

Works because systems trained on the same parallel corpus. Both halves contain the same information = both systems contain same info.

**What if no parallel data?**

Dual monolingual corpora:

**Same amount of information in each language,**  
...even if it is different information.

Adding half of a parallel sentence should also add same amount of information to both sides.

Score sentences based on monolingual entropy deltas: **are both halves equally informative** with respect to their monolingual set?

For Sinhala-English shared task: Randomly pick English data to match size of Sinhala, but same domain ratio. Use SentencePiece to force identical vocab size. Train LM, verify perplexities match on held-out (provided) parallel test set.

Corpus	Lines	Tokens	Corpus	valid	test
Sinhala Wikipedia	156k	4.7M	Sinhala, untok	1,759.8	1,565.9
Sinhala Common Crawl	5.2M	110M	English, untok	1,069.2	985.3
EN Random Wikipedia	150k	5.5M	Sinhala, tok=SP	<b>320.5</b>	<b>299.2</b>
EN Random Common Crawl	6M	123M	English, tok=SP	<b>302.5</b>	<b>292.7</b>

Cross-Entropy Delta

Measures incremental change in information of a single model. Not the same as cross-entropy difference, which is between two static models. Used as the selection criterion inside Cynical data selection (Axelrod, 2017).

How much would the perplexity of a test set change if a new sentence  $s$  was added to the training set of a language model?

$$\Delta H_{n \rightarrow n+1} = H_{n+1} - H_n$$

$$\Delta H_{n \rightarrow n+1} = \underbrace{\log \frac{W_n + w_{n+1}}{W_n}}_{Penalty} + \underbrace{\sum_{v \in V_{REPR}} \frac{C_{REPR}(v)}{W_{REPR}} \log \frac{C_n(v)}{C_n(v) + c_{n+1}(v)}}_{Gain}$$

Entropic penalty for increasing size of training set, and potential entropy improvement for adding new information to corpus.

**Use monolingual estimates of information gain instead of translation model likelihoods!**

$$|\Delta H_{E_n}(s_{E_n} | \text{MONO}_{E_n}) - \Delta H_{S_i}(s_{S_i} | \text{MONO}_{S_i})|$$

$$+ \frac{1}{2} (\Delta H_{E_n}(s_{E_n} | \text{MONO}_{E_n}) + \Delta H_{S_i}(s_{S_i} | \text{MONO}_{S_i}))$$

Results:

Successfully filtered parallel data with no parallel data! Not the best system, but it works!

Bonus: Everything runs quickly on a single laptop!

	1M SMT	1M NMT	5M SMT	5M NMT
Rank 1	4.27	6.39	4.94	4.44
DiDi	2.53	0.19	3.70	0.20
Rank 10	0.92	0.03	2.73	0.10

What happened with NMT?

User error! Length ratio feature too fancy, did not work.

System cutoffs only based on English side! Our submissions:

740k Si -> 1M En

3.6M Si -> 5M En

NMT hallucinates to produce enough output words;

BLEU score is zero. Important to get the heuristic scores right!