

# Multi-Task Learning for Improved Discriminative Training in SMT

Patrick Simianer and Stefan Riezler

Department of Computational Linguistics,  
Heidelberg University, Germany



# Learning from Big Data in SMT

- Machine learning theory and practice suggests benefits from using **expressive feature representations** and from **tuning on large training samples**.
- Discriminative training in SMT has mostly been content with tuning **small sets of dense features** on **small development data** (Och NAACL'03).
- Notable exceptions and recent success stories using **larger feature and training sets**:
  - Liang et al. ACL'06: 1.5M features, 67K parallel sentences.
  - Tillmann and Zhang ACL'06: 35M feats, 230K sents.
  - Blunsom et al. ACL'08: 7.8M feats, 100K sents.
  - Simianer, Riezler, Dyer ACL'12: 4.7M feats, 1.6M sents.
  - Flanigan, Dyer, Carbonell NAACL'13: 28.8M feats, 1M sents.

## Learning from Big Data in SMT

- Machine learning theory and practice suggests benefits from using **expressive feature representations** and from **tuning on large training samples**.
- Discriminative training in SMT has mostly been content with tuning **small sets of dense features** on **small development data** (Och NAACL'03).
- Notable exceptions and recent success stories using **larger feature and training sets**:
  - Liang et al. ACL'06: 1.5M features, 67K parallel sentences.
  - Tillmann and Zhang ACL'06: 35M feats, 230K sents.
  - Blunsom et al. ACL'08: 7.8M feats, 100K sents.
  - Simianer, Riezler, Dyer ACL'12: 4.7M feats, 1.6M sents.
  - Flanigan, Dyer, Carbonell NAACL'13: 28.8M feats, 1M sents.

## Learning from Big Data in SMT

- Machine learning theory and practice suggests benefits from using **expressive feature representations** and from **tuning on large training samples**.
- Discriminative training in SMT has mostly been content with tuning **small sets of dense features** on **small development data** (Och NAACL'03).
- Notable exceptions and recent success stories using **larger feature and training sets**:
  - Liang et al. ACL'06: 1.5M features, 67K parallel sentences.
  - Tillmann and Zhang ACL'06: 35M feats, 230K sents.
  - Blunsom et al. ACL'08: 7.8M feats, 100K sents.
  - Simianer, Riezler, Dyer ACL'12: 4.7M feats, 1.6M sents.
  - Flanigan, Dyer, Carbonell NAACL'13: 28.8M feats, 1M sents.

## Framework: Multi-Task Learning

- **Goal:** A number of statistical models need to be estimated simultaneously from data belonging to different tasks.
- **Examples:**
  - OCR of handwritten characters from different writers: Exploit commonalities on pixel- or stroke-level shared between writers.
  - LTR from search engine query logs from different countries: Some queries are country-specific ("football"), most preference rankings are shared across countries.
- **Idea:**
  - Learn a shared model that takes advantage of commonalities among tasks, without neglecting individual knowledge.
  - Problem of simultaneous learning is harder, but it also offers possibility of knowledge sharing.

## Framework: Multi-Task Learning

- **Goal:** A number of statistical models need to be estimated simultaneously from data belonging to different tasks.
- **Examples:**
  - OCR of handwritten characters from different writers: Exploit commonalities on pixel- or stroke-level shared between writers.
  - LTR from search engine query logs from different countries: Some queries are country-specific (“football”), most preference rankings are shared across countries.
- **Idea:**
  - Learn a shared model that takes advantage of commonalities among tasks, without neglecting individual knowledge.
  - Problem of simultaneous learning is harder, but it also offers possibility of knowledge sharing.

## Framework: Multi-Task Learning

- **Goal:** A number of statistical models need to be estimated simultaneously from data belonging to different tasks.
- **Examples:**
  - OCR of handwritten characters from different writers: Exploit commonalities on pixel- or stroke-level shared between writers.
  - LTR from search engine query logs from different countries: Some queries are country-specific (“football”), most preference rankings are shared across countries.
- **Idea:**
  - Learn a shared model that takes advantage of commonalities among tasks, without neglecting individual knowledge.
  - Problem of simultaneous learning is harder, but it also offers possibility of knowledge sharing.

# Multi-Task Distributed SGD for Discriminative SMT

- **Idea:** Take advantage of algorithms designed for hard problems to ease discriminative SMT on big data.
  - Distribute work,
  - learn efficiently on each example,
  - share information.
- **Method:**
  - **Distributed learning** using Hadoop/MapReduce or Sun Grid Engine.
  - **Online learning** via Stochastic Gradient Descent optimization.
  - **Feature selection** via  $\ell_1/\ell_2$  block norm regularization on features across multiple tasks.



# Multi-Task Distributed SGD for Discriminative SMT

- **Idea:** Take advantage of algorithms designed for hard problems to ease discriminative SMT on big data.
  - Distribute work,
  - learn efficiently on each example,
  - share information.
- **Method:**
  - **Distributed learning** using Hadoop/MapReduce or Sun Grid Engine.
  - **Online learning** via Stochastic Gradient Descent optimization.
  - **Feature selection** via  $\ell_1/\ell_2$  block norm regularization on features across multiple tasks.

## Related Work

- **Online learning:**
  - We deploy pairwise ranking perceptron (Shen & Joshi JMLR'05)
  - and margin perceptron (Collobert & Bengio ICML'04).
- **Distributed learning:**
  - Without feature selection, our algorithm reduces to Iterative Mixing (McDonald et al. NAACL'10),
  - which itself is related to Bagging (Breiman JMLR'96) if shards are treated as random samples.

## Related Work

- **Online learning:**
  - We deploy pairwise ranking perceptron (Shen & Joshi JMLR'05)
  - and margin perceptron (Collobert & Bengio ICML'04).
- **Distributed learning:**
  - Without feature selection, our algorithm reduces to Iterative Mixing (McDonald et al. NAACL'10),
  - which itself is related to Bagging (Breiman JMLR'96) if shards are treated as random samples.

## Related Work

- $\ell_1/\ell_2$  **regularization**:
  - Related to group-Lasso approaches which use mixed norms (Yuan & Lin JRSS'06), hierarchical norms (Zhao et al. Annals Stats'09), structured norms (Martins et al. EMNLP'11).
  - Difference: Norms and proximity operators are applied to *groups* of features in *single* regression or classification task – multi-task learning groups features orthogonally by tasks.
  - Closest relation to Obozinski et al. StatComput'10: Our algorithm is weight-based backward feature elimination variant of their gradient-based forward feature selection algorithm.

## OL Framework: Pairwise Ranking Perceptron

- Preference pairs  $\mathbf{x}_j = (\mathbf{x}_j^{(1)}, \mathbf{x}_j^{(2)})$  where  $\mathbf{x}_j^{(1)}$  is ordered above  $\mathbf{x}_j^{(2)}$  w.r.t. sentence-wise BLEU (Nakov et al. COLING'12).
- Hinge loss-type objective

$$l_j(\mathbf{w}) = (-\langle \mathbf{w}, \bar{\mathbf{x}}_j \rangle)_+$$

where  $\bar{\mathbf{x}}_j = \mathbf{x}_j^{(1)} - \mathbf{x}_j^{(2)}$ ,  $(a)_+ = \max(0, a)$ ,  $\mathbf{w} \in \mathbb{R}^D$  is a weight vector, and  $\langle \cdot, \cdot \rangle$  denotes the standard vector dot product.

- **Ranking perceptron** by stochastic subgradient descent:

$$\nabla l_j(\mathbf{w}) = \begin{cases} -\bar{\mathbf{x}}_j & \text{if } \langle \mathbf{w}, \bar{\mathbf{x}}_j \rangle \leq 0, \\ 0 & \text{else.} \end{cases}$$

## OL Framework: Pairwise Ranking Perceptron

- Preference pairs  $\mathbf{x}_j = (\mathbf{x}_j^{(1)}, \mathbf{x}_j^{(2)})$  where  $\mathbf{x}_j^{(1)}$  is ordered above  $\mathbf{x}_j^{(2)}$  w.r.t. sentence-wise BLEU (Nakov et al. COLING'12).
- Hinge loss-type objective

$$l_j(\mathbf{w}) = (-\langle \mathbf{w}, \bar{\mathbf{x}}_j \rangle)_+$$

where  $\bar{\mathbf{x}}_j = \mathbf{x}_j^{(1)} - \mathbf{x}_j^{(2)}$ ,  $(a)_+ = \max(0, a)$ ,  $\mathbf{w} \in \mathbb{R}^D$  is a weight vector, and  $\langle \cdot, \cdot \rangle$  denotes the standard vector dot product.

- Ranking perceptron by stochastic subgradient descent:

$$\nabla l_j(\mathbf{w}) = \begin{cases} -\bar{\mathbf{x}}_j & \text{if } \langle \mathbf{w}, \bar{\mathbf{x}}_j \rangle \leq 0, \\ 0 & \text{else.} \end{cases}$$

## OL Framework: Pairwise Ranking Perceptron

- Preference pairs  $\mathbf{x}_j = (\mathbf{x}_j^{(1)}, \mathbf{x}_j^{(2)})$  where  $\mathbf{x}_j^{(1)}$  is ordered above  $\mathbf{x}_j^{(2)}$  w.r.t. sentence-wise BLEU (Nakov et al. COLING'12).
- Hinge loss-type objective

$$l_j(\mathbf{w}) = (-\langle \mathbf{w}, \bar{\mathbf{x}}_j \rangle)_+$$

where  $\bar{\mathbf{x}}_j = \mathbf{x}_j^{(1)} - \mathbf{x}_j^{(2)}$ ,  $(a)_+ = \max(0, a)$ ,  $\mathbf{w} \in \mathbb{R}^D$  is a weight vector, and  $\langle \cdot, \cdot \rangle$  denotes the standard vector dot product.

- **Ranking perceptron** by stochastic subgradient descent:

$$\nabla l_j(\mathbf{w}) = \begin{cases} -\bar{\mathbf{x}}_j & \text{if } \langle \mathbf{w}, \bar{\mathbf{x}}_j \rangle \leq 0, \\ 0 & \text{else.} \end{cases}$$

## OL framework: Margin Perceptron

- Hinge loss-type objective

$$l_j(\mathbf{w}) = (1 - \langle \mathbf{w}, \bar{\mathbf{x}}_j \rangle)_+$$

- Stochastic subgradient descent:

$$\nabla l_j(\mathbf{w}) = \begin{cases} -\bar{\mathbf{x}}_j & \text{if } \langle \mathbf{w}, \bar{\mathbf{x}}_j \rangle < 1, \\ 0 & \text{else.} \end{cases}$$

- Margin term controls capacity, but results in more updates.
- Collobert & Bengio (ICML'04) argue that this justifies not using an explicit regularization (as for example in an SGD version of the SVM (Shalev-Shwartz et al. ICML'07)).



## OL framework: Margin Perceptron

- Hinge loss-type objective

$$l_j(\mathbf{w}) = (1 - \langle \mathbf{w}, \bar{\mathbf{x}}_j \rangle)_+$$

- Stochastic subgradient descent:

$$\nabla l_j(\mathbf{w}) = \begin{cases} -\bar{\mathbf{x}}_j & \text{if } \langle \mathbf{w}, \bar{\mathbf{x}}_j \rangle < 1, \\ 0 & \text{else.} \end{cases}$$

- Margin term controls capacity, but results in more updates.
- Collobert & Bengio (ICML'04) argue that this justifies not using an explicit regularization (as for example in an SGD version of the SVM (Shalev-Shwartz et al. ICML'07)).

## OL framework: Margin Perceptron

- Hinge loss-type objective

$$l_j(\mathbf{w}) = (1 - \langle \mathbf{w}, \bar{\mathbf{x}}_j \rangle)_+$$

- Stochastic subgradient descent:

$$\nabla l_j(\mathbf{w}) = \begin{cases} -\bar{\mathbf{x}}_j & \text{if } \langle \mathbf{w}, \bar{\mathbf{x}}_j \rangle < 1, \\ 0 & \text{else.} \end{cases}$$

- Margin term controls capacity, but results in more updates.
- Collobert & Bengio (ICML'04) argue that this justifies not using an explicit regularization (as for example in an SGD version of the SVM (Shalev-Shwartz et al. ICML'07)).

## MTL Framework: $\ell_1/\ell_2$ Block Norm Regularization

- Data points  $\{(x_{zn}, y_{zn}), z = 1, \dots, Z, n = 1, \dots, N_z\}$ , sampled from  $P_z$  on  $X \times Y$  ( $z = \text{task}; n = \text{data point}$ ).
- Objective:

$$\min_{\mathbf{W}} \sum_{z,n} l_n(\mathbf{w}_z) + \lambda \|\mathbf{W}\|_{1,2}$$

- where  $\mathbf{W} = (\mathbf{w}_z^d)_{z,d}$  is a  $Z$ -by- $D$  matrix  $\mathbf{W} = (\mathbf{w}_z^d)_{z,d}$  of  $D$ -dimensional row vectors  $\mathbf{w}_z$  and  $Z$ -dimensional column vectors  $\mathbf{w}^d$  of weights associated with feature  $d$  across tasks.
- Weighted  $\ell_1/\ell_2$  norm:

$$\lambda \|\mathbf{W}\|_{1,2} = \lambda \sum_{d=1}^D \|\mathbf{w}^d\|_2$$

- Each  $\ell_2$  norm of a weight column  $\mathbf{w}^d$  represents the relevance of the corresponding feature across tasks.

## MTL Framework: $\ell_1/\ell_2$ Block Norm Regularization

- Data points  $\{(x_{zn}, y_{zn}), z = 1, \dots, Z, n = 1, \dots, N_z\}$ , sampled from  $P_z$  on  $X \times Y$  ( $z = \text{task}$ ;  $n = \text{data point}$ ).
- Objective:

$$\min_{\mathbf{W}} \sum_{z,n} l_n(\mathbf{w}_z) + \lambda \|\mathbf{W}\|_{1,2}$$

- where  $\mathbf{W} = (\mathbf{w}_z^d)_{z,d}$  is a  $Z$ -by- $D$  matrix  $\mathbf{W} = (\mathbf{w}_z^d)_{z,d}$  of  $D$ -dimensional row vectors  $\mathbf{w}_z$  and  $Z$ -dimensional column vectors  $\mathbf{w}^d$  of weights associated with feature  $d$  across tasks.
- Weighted  $\ell_1/\ell_2$  norm:

$$\lambda \|\mathbf{W}\|_{1,2} = \lambda \sum_{d=1}^D \|\mathbf{w}^d\|_2$$

- Each  $\ell_2$  norm of a weight column  $\mathbf{w}^d$  represents the relevance of the corresponding feature across tasks.

## MTL Framework: $\ell_1/\ell_2$ Block Norm Regularization

- Data points  $\{(x_{zn}, y_{zn}), z = 1, \dots, Z, n = 1, \dots, N_z\}$ , sampled from  $P_z$  on  $X \times Y$  ( $z = \text{task}$ ;  $n = \text{data point}$ ).
- Objective:

$$\min_{\mathbf{W}} \sum_{z,n} l_n(\mathbf{w}_z) + \lambda \|\mathbf{W}\|_{1,2}$$

- where  $\mathbf{W} = (\mathbf{w}_z^d)_{z,d}$  is a  $Z$ -by- $D$  matrix  $\mathbf{W} = (\mathbf{w}_z^d)_{z,d}$  of  $D$ -dimensional row vectors  $\mathbf{w}_z$  and  $Z$ -dimensional column vectors  $\mathbf{w}^d$  of weights associated with feature  $d$  across tasks.
- Weighted  $\ell_1/\ell_2$  norm:

$$\lambda \|\mathbf{W}\|_{1,2} = \lambda \sum_{d=1}^D \|\mathbf{w}^d\|_2$$

- Each  $\ell_2$  norm of a weight column  $\mathbf{w}^d$  represents the relevance of the corresponding feature across tasks.

# $\ell_1/\ell_2$ Regularization Explained

	$\mathbf{w}^1$	$\mathbf{w}^2$	$\mathbf{w}^3$	$\mathbf{w}^4$	$\mathbf{w}^5$		$\mathbf{w}^1$	$\mathbf{w}^2$	$\mathbf{w}^3$	$\mathbf{w}^4$	$\mathbf{w}^5$		
$\mathbf{w}_{z_1}$	[	6	4	0	0	0	[	6	4	0	0	0	]
$\mathbf{w}_{z_2}$	[	0	0	3	0	0	[	3	0	0	0	0	]
$\mathbf{w}_{z_3}$	[	0	0	0	2	3	[	2	3	0	0	0	]
column $\ell_2$ norm:		6	4	3	2	3		7	5	0	0	0	
$\ell_1$ sum:						$\Rightarrow$						$\Rightarrow$	12

- $\ell_1$  sum of  $\ell_2$  norms encourages several feature columns  $\mathbf{w}^d$  to be  $\mathbf{0}$  and others to have high weights across tasks.
- **Algorithm idea:**

- Contribution to loss reduction must outweigh regularizer penalty in order to activate feature by non-zero weight.
- Weight-based feature elimination criterion:

$$\text{If } \|\mathbf{w}^d\|_2 \leq \lambda, \text{ set } \mathbf{W}[z][d] = 0, \forall z.$$

- Implementation by threshold on  $K$  features or by threshold  $\lambda$ .

## $\ell_1/\ell_2$ Regularization Explained

	$\mathbf{w}^1$	$\mathbf{w}^2$	$\mathbf{w}^3$	$\mathbf{w}^4$	$\mathbf{w}^5$		$\mathbf{w}^1$	$\mathbf{w}^2$	$\mathbf{w}^3$	$\mathbf{w}^4$	$\mathbf{w}^5$		
$\mathbf{w}_{z_1}$	[	6	4	0	0	0	[	6	4	0	0	0	]
$\mathbf{w}_{z_2}$	[	0	0	3	0	0	[	3	0	0	0	0	]
$\mathbf{w}_{z_3}$	[	0	0	0	2	3	[	2	3	0	0	0	]
column $\ell_2$ norm:		6	4	3	2	3		7	5	0	0	0	
$\ell_1$ sum:						$\Rightarrow$						$\Rightarrow$	12

- $\ell_1$  sum of  $\ell_2$  norms encourages several feature columns  $\mathbf{w}^d$  to be  $\mathbf{0}$  and others to have high weights across tasks.
- **Algorithm idea:**
  - Contribution to loss reduction must outweigh regularizer penalty in order to activate feature by non-zero weight.
  - Weight-based feature elimination criterion:

$$\text{If } \|\mathbf{w}^d\|_2 \leq \lambda, \text{ set } \mathbf{W}[z][d] = 0, \forall z.$$

- Implementation by threshold on  $K$  features or by threshold  $\lambda$ .

# Implementation as Feature Selection Algorithm

---

## Algorithm 1 Multi-task Distributed SGD

---

Get data for  $Z$  tasks, each including  $S$  sentences;  
 distribute to machines.

Initialize  $\mathbf{v} \leftarrow \mathbf{0}$ .

```

for epochs  $t \leftarrow 0 \dots T - 1$ : do
  for all tasks  $z \in \{1 \dots Z\}$ : parallel do
     $\mathbf{w}_{z,t,0,0} \leftarrow \mathbf{v}$ 
    for all sentences  $i \in \{0 \dots S - 1\}$ : do
      Decode  $i^{\text{th}}$  input with  $\mathbf{w}_{z,t,i,0}$ .
      for all pairs  $j \in \{0 \dots P - 1\}$ : do
         $\mathbf{w}_{z,t,i,j+1} \leftarrow \mathbf{w}_{z,t,i,j} - \eta \nabla_{\mathbf{w}_{z,t,i,j}}$ 
      end for
       $\mathbf{w}_{z,t,i+1,0} \leftarrow \mathbf{w}_{z,t,i,P}$ 
    end for
  end for
  Stack weights  $\mathbf{W} \leftarrow [\mathbf{w}_{1,t,S,0} \mid \dots \mid \mathbf{w}_{Z,t,S,0}]^T$ 
  Select top  $K$  feature columns of  $\mathbf{W}$  by  $\ell_2$  norm
  for  $k \leftarrow 1 \dots K$  do
     $\mathbf{v}[k] = \frac{1}{Z} \sum_{z=1}^Z \mathbf{W}[z][k]$ 
  end for
end for
return  $\mathbf{v}$ 
  
```

---



# Experiments: Random vs. Natural Tasks

- **Research Question:**

- As shown in earlier work (Simianer, Riezler, Dyer ACL'12), multi-task learning can be used as general regularization technique on **random shards**.
- Can multi-task learning benefit from **natural task structure** in the data, where shared and individual knowledge is properly balanced?

# Experiments: Random vs. Natural Tasks

- **Research Question:**
  - As shown in earlier work (Simianer, Riezler, Dyer ACL'12), multi-task learning can be used as general regularization technique on **random shards**.
  - Can multi-task learning benefit from **natural task structure** in the data, where shared and individual knowledge is properly balanced?

# Data

---

A	Human Necessities
B	Performing Operations, Transporting
C	Chemistry, Metallurgy
D	Textiles, Paper
E	Fixed Constructions
F	Mechanical Engineering, Lighting, Heating, Weapons
G	Physics
H	Electricity

---

- International Patent Classification (IPC) categorizes patents hierarchically into eight sections, 120 classes, 600 subclasses, down to 70,000 subgroups at the leaf level.
- Typically, a patent belongs to more than one section, with one section chosen as main classification.
- Eight top classes/sections used to define **natural tasks**.

## SMT Setup

- (1)  $X \rightarrow X_1$  hat  $X_2$  versprochen;  $X_1$  promised  $X_2$
  - (2)  $X \rightarrow X_1$  hat mir  $X_2$  versprochen;  
 $X_1$  promised me  $X_2$
  - (3)  $X \rightarrow X_1$  versprach  $X_2$ ;  $X_1$  promised  $X_2$
- Hierarchical phrase-based translation (Chiang CL'07), formalizes translation rules as productions of synchronous context-free grammar (SCFG).
  - Features in discriminative training:
    - **Rule identifiers** for SCFG productions  
Examples: rule (1), (2) and (3)
    - **Rule n-gram** features in source and target  
Examples: "X hat", "hat X", "X versprochen"
    - **Rule shape** features  
Examples: (NT, term\*, NT, term\*; NT, term\*, NT) for (1), (2);  
(NT, term\*, NT; NT, term\*, NT) for rule (3).

## SMT Setup

(1)  $X \rightarrow X_1 \text{ hat } X_2 \text{ versprochen}; X_1 \text{ promised } X_2$

(2)  $X \rightarrow X_1 \text{ hat mir } X_2 \text{ versprochen};$

$X_1 \text{ promised me } X_2$

(3)  $X \rightarrow X_1 \text{ versprach } X_2; X_1 \text{ promised } X_2$

- Hierarchical phrase-based translation (Chiang CL'07), formalizes translation rules as productions of synchronous context-free grammar (SCFG).
- Features in discriminative training:
  - **Rule identifiers** for SCFG productions  
Examples: rule (1), (2) and (3)
  - **Rule n-gram** features in source and target  
Examples: “X hat”, “hat X”, “X versprochen”
  - **Rule shape** features  
Examples: (NT, term\*, NT, term\*; NT, term\*, NT) for (1), (2);  
(NT, term\*, NT; NT, term\*, NT) for rule (3).

## MERT Baseline w/ 12 Dense Features

	single-task tuning		
	indep. <sup>0</sup>	pooled <sup>1</sup>	pooled-cat <sup>2</sup>
pooled test	–	51.18	51.22
A	54.92	<sup>02</sup> 55.27	<sup>0</sup> 55.17
B	51.53	51.48	<sup>01</sup> 51.69
C	<sup>12</sup> 56.31	<sup>2</sup> 55.90	55.74
D	49.94	<sup>0</sup> 50.33	<sup>0</sup> 50.26
E	<sup>1</sup> 49.19	48.97	<sup>1</sup> 49.13
F	<sup>12</sup> 51.26	51.02	51.12
G	<sup>1</sup> 49.61	49.44	49.55
H	49.38	49.50	<sup>01</sup> 49.67
<b>average test</b>	<b>51.52</b>	<b>51.49</b>	<b>51.54</b>

- Neither tuning on *pooled* or *pooled-cat* improves over *indep.*.
- $x \in \{0,1,2\}$  BLEU denotes statistical significance of pairwise test.
- Tuning was repeated 3 times and BLEU scores averaged.

## Single-Task Perceptron w/ $\ell_1$ Regularization

	single-task tuning		
	indep. <sup>0</sup>	pooled <sup>1</sup>	pooled-cat <sup>2</sup>
pooled test	–	50.75	<sup>1</sup> 52.08
A	<sup>1</sup> 55.11	54.32	<sup>01</sup> 55.94
B	<sup>1</sup> 52.61	50.84	<sup>1</sup> 52.57
C	56.18	56.11	<sup>01</sup> 56.75
D	<sup>1</sup> 50.68	49.48	<sup>01</sup> 51.22
E	<sup>1</sup> 50.27	48.69	<sup>1</sup> 50.01
F	<sup>1</sup> 51.68	50.71	<sup>1</sup> 51.95
G	<sup>1</sup> 49.90	49.06	<sup>01</sup> 50.51
H	<sup>1</sup> 50.48	49.16	<sup>1</sup> 50.53
<b>average test</b>	<b>52.11</b>	<b>51.05</b>	<b>52.44</b>
model size	430,092.5	457,428	1,574,259

- Improvements over MERT, mostly on *pooled-cat* tuning set.
- 1.5M features make serial tuning on *pooled-cat* infeasible.
- Overfitting effect on small *pooled* data.

## Single- and Multi-Task Perceptron

	single-task tuning			multi-task tuning		
	indep. <sup>0</sup>	pooled <sup>1</sup>	pooled-cat <sup>2</sup>	IPC <sup>3</sup>	sharding <sup>4</sup>	resharding <sup>5</sup>
pooled test	–	51.33	1 51.77	<sup>12</sup> 52.56	<sup>12</sup> 52.54	<sup>12</sup> 52.60
A	54.79	54.76	<sup>01</sup> 55.31	<sup>012</sup> 56.35	<sup>012</sup> 56.22	<sup>012</sup> 56.21
B	<sup>12</sup> 52.45	51.30	<sup>1</sup> 52.19	<sup>012</sup> 52.78	<sup>0123</sup> 52.98	<sup>012</sup> 52.96
C	<sup>2</sup> 56.62	56.65	<sup>1</sup> 56.12	<sup>01245</sup> 57.76	<sup>012</sup> 57.30	<sup>012</sup> 57.44
D	<sup>1</sup> 50.75	49.88	<sup>1</sup> 50.63	<sup>01245</sup> 51.54	<sup>012</sup> 51.33	<sup>012</sup> 51.20
E	<sup>1</sup> 49.70	49.23	<sup>01</sup> 49.92	<sup>012</sup> 50.51	<sup>012</sup> 50.52	<sup>012</sup> 50.38
F	<sup>1</sup> 51.60	51.09	<sup>1</sup> 51.71	<sup>012</sup> 52.28	<sup>012</sup> 52.43	<sup>012</sup> 52.32
G	<sup>1</sup> 49.50	49.06	<sup>01</sup> 49.97	<sup>012</sup> 50.84	<sup>012</sup> 50.88	<sup>012</sup> 50.74
H	<sup>1</sup> 49.77	49.50	<sup>01</sup> 50.64	<sup>012</sup> 51.16	<sup>012</sup> 51.07	<sup>012</sup> 51.10
<b>average test</b>	<b>51.90</b>	<b>51.42</b>	<b>52.06</b>	<b>52.90</b>	<b>52.84</b>	<b>52.79</b>
model size	366,869.4	448,359	1,478,049	100,000	100,000	100,000

- Multi-task tuning improves BLEU over all single-task runs.
- Also more efficient due to iterative feature selection.
- Difference between natural and random tasks inconclusive.



# Single- and Multi-Task Margin Perceptron

	single-task tuning			multi-task tuning		
	indep. <sup>0</sup>	pooled <sup>1</sup>	pooled-cat <sup>2</sup>	IPC <sup>3</sup>	sharding <sup>4</sup>	resharding <sup>5</sup>
pooled test	–	51.33	<sup>1</sup> 52.58	<sup>12</sup> 52.98	<sup>12</sup> 52.95	<sup>12</sup> 52.99
A	<sup>1</sup> 56.09	55.33	<sup>1</sup> 55.92	<sup>01245</sup> 56.78	<sup>012</sup> 56.62	<sup>012</sup> 56.53
B	<sup>1</sup> 52.45	51.59	<sup>1</sup> 52.44	<sup>012</sup> 53.31	<sup>012</sup> 53.35	<sup>012</sup> 53.21
C	<sup>1</sup> 57.20	56.85	<sup>01</sup> 57.54	<sup>01</sup> 57.46	<sup>1</sup> 57.42	<sup>1</sup> 57.43
D	<sup>1</sup> 50.51	50.18	<sup>01</sup> 51.38	<sup>01245</sup> 52.14	<sup>0125</sup> 51.82	<sup>012</sup> 51.66
E	<sup>1</sup> 50.27	49.36	<sup>01</sup> 50.72	<sup>0124</sup> 51.13	<sup>012</sup> 50.89	<sup>012</sup> 51.02
F	<sup>1</sup> 52.06	51.20	<sup>01</sup> 52.61	<sup>01245</sup> 53.07	<sup>012</sup> 52.80	<sup>012</sup> 52.87
G	<sup>1</sup> 50.00	49.58	<sup>01</sup> 50.90	<sup>01245</sup> 51.36	<sup>012</sup> 51.19	<sup>012</sup> 51.11
H	<sup>1</sup> 50.57	49.80	<sup>01</sup> 51.32	<sup>012</sup> 51.57	<sup>012</sup> 51.62	<sup>01</sup> 51.47
<b>average test</b>	<b>52.39</b>	<b>51.74</b>	<b>52.85</b>	<b>53.35</b>	<b>53.21</b>	<b>53.16</b>
model size	423,731.5	484,483	1,697,398	100,000	100,000	100,000

- Single-task runs beat standard perceptron w/ and w/o  $\ell_1$ .
- Regularization by margin and multi-task learning adds up.
- Best result is nearly 2 BLEU points better than MERT.

# Conclusion

- Multi-task learning for SMT is **efficient** due to online learning, parallelization and feature selection,
- but also **effective** in terms of BLEU improvements over single-task learning.
- Multi-task distributed learning is **easy to implement as wrapper** around perceptron.

## Future Work: Task Adaption

- *Natural* tasks are slightly advantageous over *random* tasks.
- Goal: Adapt task definition to SMT problem.
  - Explore various similarity metrics on IPC subclasses,
  - jointly optimize task partitioning and SMT performance.

## Future Work: Adaptive Regularization

---

### Algorithm 2 Path-Following Multi-task Distributed SGD

---

Get data for  $Z$  tasks, each including  $S$  sentences; distribute to machines.

Initialize  $\mathbf{v} \leftarrow \mathbf{0}$ ;  $\lambda_0, \lambda_{\min}, \epsilon$ .

for epochs  $t \leftarrow 0 \dots T - 1$ : do

  for all tasks  $z \in \{1 \dots Z\}$ : parallel do

    Perform task-specific learning

  end for

  Stack weights  $\mathbf{W} \leftarrow [\mathbf{w}_{1,t,S,0} | \dots | \mathbf{w}_{Z,t,S,0}]^T$

  for feature columns  $d \in \{1 \dots D\}$  in  $\mathbf{W}$ : do

    if  $\|\mathbf{w}^d\|_2 \leq \lambda_t$  then

$\mathbf{v}[d] = 0$

    else

$$\mathbf{v}[d] = \frac{1}{Z} \sum_{z=1}^Z \mathbf{W}[z][d]$$

    end if

  end for

  Set  $\lambda_{t+1} = \min\left\{\lambda_t, \frac{\sum_{z,i,j} (l_{z,i,j}(\mathbf{v}_{t-1}) - l_{z,i,j}(\mathbf{v}_t))}{\epsilon}\right\}$

  if  $\lambda_{t+1} < \lambda_{\min}$  then

    break

  end if

end for

return  $\mathbf{v}$

---

# Thanks for your attention!

`dtrain codebase is part of cdec:`  
`https://github.com/redpony/cdec.`