



# Multitask Spectral Learning of Weighted Finite Automata

Guillaume Rabusseau<sup>1</sup> Borja Balle<sup>2</sup> Joelle Pineau<sup>1,3</sup>

<sup>1</sup>McGill University <sup>2</sup>Amazon Research Cambridge <sup>3</sup>CIFAR



## Overview

- Sequence data is ubiquitous in computer science and machine learning.



- Weighted Finite Automata (WFA) can model functions on sequences.
- We propose a **spectral multitask learning** algorithm for WFAs:
  - extends the spectral learning algorithm for WFAs [1]
  - relies on the novel model of **vector-valued WFA**.

## Multitask Learning

- Common task in machine learning: estimate a function  $f: \mathcal{X} \rightarrow \mathcal{Y}$  from a training sample  $\{(x_i, y_i)\}_{i=1}^N$  where each  $y_i \simeq f(x_i)$ .
- In **multitask learning**, the learner is given several learning tasks  $f_1, \dots, f_m$ .
- Jointly learning **related tasks**  $f_1, \dots, f_m$  can lead to better performances.
- This work: multitask learning when  $\mathcal{X}$  consists of *sequence data*.

## Weighted Finite Automata

- A **weighted finite automaton** (WFA) is a tuple  $A = (\alpha, \{\mathbf{A}^\sigma\}_{\sigma \in \Sigma}, \omega)$  and computes a function  $f_A: \Sigma^* \rightarrow \mathbb{R}$  defined for each word  $x = x_1 x_2 \dots x_k \in \Sigma^*$  by

$$f_A(x_1 x_2 \dots x_k) = \alpha^\top \mathbf{A}^{x_1} \mathbf{A}^{x_2} \dots \mathbf{A}^{x_k} \omega = \alpha^\top \mathbf{A}^x \omega.$$

- The **number of states** of  $A$  is the size  $n$  of the matrices  $\mathbf{A}^\sigma$  and  $A$  is **minimal** if any WFA  $B$  such that  $f_A = f_B$  has at least  $n$  states, in which case  $n$  is the **rank** of the function  $f$ .

## Spectral Learning of WFAs

- Hankel matrix**  $\mathbf{H}_f \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$  associated with a function  $f: \Sigma^* \rightarrow \mathbb{R}$
- $$(\mathbf{H}_f)_{u,v} = f(uv) \text{ for all } u, v \in \Sigma^*.$$

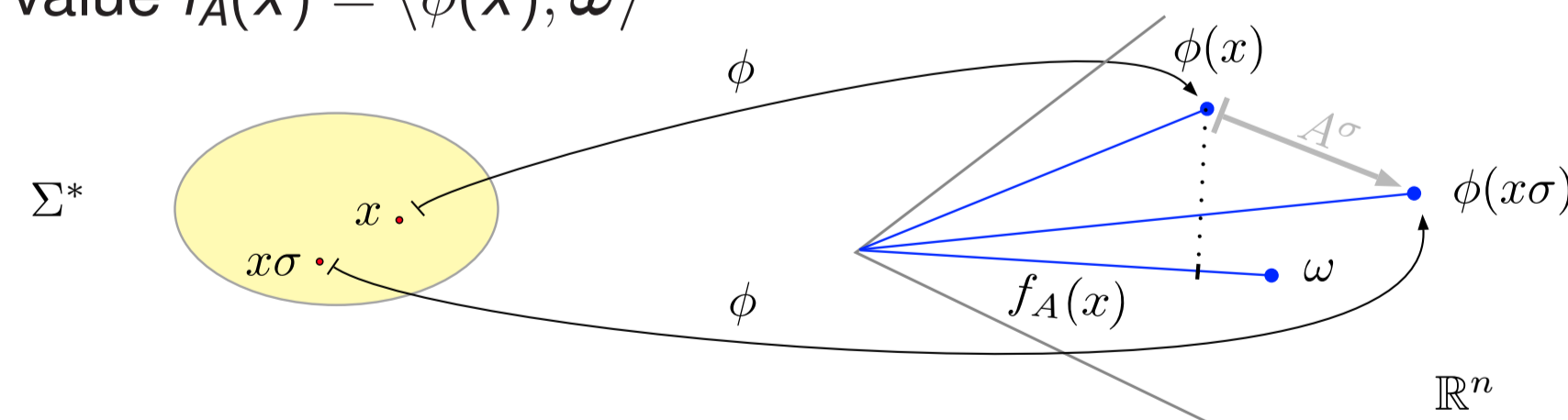
**Theorem** [3, 4] For any function  $f: \Sigma^* \rightarrow \mathbb{R}$ ,  $\text{rank}(f) = \text{rank}(\mathbf{H}_f)$ .

- Spectral learning** of WFAs (in a nutshell) [1, Lemma 4.1].

- Let  $\mathbf{H}_f = \mathbf{P}\mathbf{S}$  with  $\mathbf{P}, \mathbf{S}^\top \in \mathbb{R}^{\Sigma^* \times n}$  where  $n = \text{rank}(f)$
- For each  $\sigma \in \Sigma$ , let  $\mathbf{H}_f^\sigma \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$  be defined by  $(\mathbf{H}_f^\sigma)_{u,v} = f(u\sigma v)$  for all  $u, v \in \Sigma^*$ .
- WFA  $(\alpha, \{\mathbf{A}^\sigma\}_{\sigma \in \Sigma}, \omega)$  with  $\alpha^\top = \mathbf{P}_{\lambda,:}$ ,  $\omega = \mathbf{S}_{:, \lambda}$ , and  $\mathbf{A}^\sigma = \mathbf{P}^\dagger \mathbf{H}_f^\sigma \mathbf{S}^\dagger$  is a **minimal WFA for  $f$** .

## WFAs as Linear Models in a Feature Space

- Computation of a WFA  $A$  on  $x \in \Sigma^*$ :
  - map  $x$  to feature vector  $\phi(x) = \alpha^\top \mathbf{A}^x$  through a **compositional feature map**  $\phi: \Sigma^* \rightarrow \mathbb{R}^n$
  - compute final value  $f_A(x) = \langle \phi(x), \omega \rangle$



- $\phi$  is **compositional**:  $\phi(x\sigma)^\top = \phi(x)^\top \mathbf{A}^\sigma$ .
- $\phi$  is **minimal** if  $V = \text{span}(\{\phi(x)\}_{x \in \Sigma^*}) \subset \mathbb{R}^n$  is of dimension  $n$ .
- $\phi: x \mapsto \alpha^\top \mathbf{A}^x$  is minimal if and only if  $(\alpha, \{\mathbf{A}^\sigma\}_{\sigma \in \Sigma}, \omega)$  is minimal.

## A Notion of Relatedness between Functions on Sequences

**Relatedness between WFAs**: to which extent two WFAs can share a joint feature map  $\phi$ .

- Let  $f_1, f_2: \Sigma^* \rightarrow \mathbb{R}$  of rank  $n_1$  and  $n_2$ . with feature maps  $\phi_1: \Sigma^* \rightarrow \mathbb{R}^{n_1}$  and  $\phi_2: \Sigma^* \rightarrow \mathbb{R}^{n_2}$ .
- $\phi = \phi_1 \oplus \phi_2: \Sigma^* \rightarrow \mathbb{R}^{n_1+n_2}$  is a joint feature map for  $f_1$  and  $f_2$ :

$$f_1(x) = \langle \phi(x), \omega_1 \oplus \mathbf{0} \rangle \text{ and } f_2(x) = \langle \phi(x), \mathbf{0} \oplus \omega_2 \rangle$$

but it may not be minimal.

→ there may exist another feature map of dimension  $n < n_1 + n_2$ .

- The smaller  $n$  is, the more related  $f_1$  and  $f_2$  are.

## Vector-Valued WFA

- A  $d$ -dimensional **vector-valued weighted finite automaton** (vv-WFA) with  $n$  states is a tuple  $A = (\alpha, \{\mathbf{A}^\sigma\}_{\sigma \in \Sigma}, \Omega)$  where
  - $\alpha \in \mathbb{R}^n$  is the **initial weights vector**
  - $\Omega \in \mathbb{R}^{n \times d}$  is the **matrix of final weights**
  - $\mathbf{A}^\sigma \in \mathbb{R}^{n \times n}$  is the **transition matrix** for each  $\sigma \in \Sigma$ .
- A vv-WFA computes a function  $\vec{f}_A: \Sigma^* \rightarrow \mathbb{R}^d$  defined for each word  $x = x_1 x_2 \dots x_k \in \Sigma^*$  by

$$\vec{f}_A(x_1 x_2 \dots x_k) = \alpha^\top \mathbf{A}^{x_1} \mathbf{A}^{x_2} \dots \mathbf{A}^{x_k} \Omega = \alpha^\top \mathbf{A}^x \Omega.$$

- Rank of  $\vec{f} = [f_1, f_2]: \Sigma^* \rightarrow \mathbb{R}^2$  equal dimension of a minimal joint feature map for  $f_1$  and  $f_2$ .
- $\max\{\text{rank}(f_1), \text{rank}(f_2)\} \leq \text{rank}([f_1, f_2]) \leq \text{rank}(f_1) + \text{rank}(f_2)$ .

**Example**

$$\begin{cases} f_1(x) = 0.5|x|_a + 0.5|x|_b & \text{rank } f_2 = 4 = \text{rank}([f_2, f_3]) \\ f_2(x) = 0.3|x|_b - 0.6|x|_c & \text{rank}([f_1, f_3]) = 6 = \text{rank}(f_1) + \text{rank}(f_3) \\ f_3(x) = |x|_c & \text{rank}(f_1) = \text{rank}(f_2) < \text{rank}([f_1, f_2]) < \text{rank}(f_1) + \text{rank}(f_2) \end{cases}$$

## Spectral Learning of vv-WFAs

- Hankel tensor**  $\mathcal{H} \in \mathbb{R}^{\Sigma^* \times d \times \Sigma^*}$  associated with a function  $\vec{f}: \Sigma^* \rightarrow \mathbb{R}^d$
- $$\mathcal{H}_{u,v,v} = \vec{f}(uv) \text{ for all } u, v \in \Sigma^*.$$

**Theorem** [Vector-Valued Fliess Theorem] For any  $\vec{f}: \Sigma^* \rightarrow \mathbb{R}^d$ ,  $\text{rank}(\vec{f}) = \text{rank}(\mathcal{H}_{(1)})$ , where  $\mathcal{H}_{(1)} = [\mathcal{H}_{:,1,:}, \mathcal{H}_{:,2,:}, \dots, \mathcal{H}_{:,d,:}]$  is the flattening of the Hankel tensor.

- Spectral learning** of vv-WFAs. A vv-WFA computing  $\vec{f}$  can be recovered from any rank  $n$  factorization of  $\mathcal{H}_{(1)}$ :

- Let  $\mathcal{H}_{(1)} = \mathbf{P}\mathbf{S}$  with  $\mathbf{P} \in \mathbb{R}^{\Sigma^* \times n}$  and  $\mathbf{S} \in \mathbb{R}^{n \times d \times \Sigma^*}$ .
- For each  $\sigma \in \Sigma$ , let  $\mathcal{H}^\sigma \in \mathbb{R}^{\Sigma^* \times d \times \Sigma^*}$  be defined by  $\mathcal{H}_{u,v,v}^\sigma = \vec{f}(u\sigma v)$  for all  $u, v \in \Sigma^*$ .
- The vv-WFA  $A = (\alpha, \{\mathbf{A}^\sigma\}_{\sigma \in \Sigma}, \Omega)$  where  $\alpha^\top = \mathbf{P}_{\lambda,:}$ ,  $\Omega = \mathbf{S}_{:, \lambda}$ , and  $\mathbf{A}^\sigma = \mathbf{P}^\dagger \mathcal{H}_{(1)}^\sigma (\mathbf{S}_{(1)})^\dagger$  is a minimal vv-WFA for  $\vec{f}$ .

## Multitask Learning of WFAs

- Let  $f_1, \dots, f_m$  be related functions defined on  $\Sigma^*$ .
- Learning  $\vec{f} = [f_1, \dots, f_m]$  as a vv-WFA **enforces discovering a shared feature map** between tasks.

## Algorithm 1

- Compute the rank  $R$  truncated SVD  $\hat{\mathcal{H}}_{(1)} \simeq \mathbf{U}\mathbf{D}\mathbf{V}^\top$ .
- Let  $A = (\alpha, \{\mathbf{A}^\sigma\}_{\sigma \in \Sigma}, \Omega)$  be the vv-WFA defined by  $\alpha^\top = \mathbf{U}_{\lambda,:}$ ,  $\Omega = \mathbf{U}^\top (\hat{\mathcal{H}}_{(1)})_{:, \lambda}$  and  $\mathbf{A}^\sigma = \mathbf{U}^\top \hat{\mathcal{H}}_{(1)}^\sigma (\hat{\mathcal{H}}_{(1)})^\dagger \mathbf{U}$  for each  $\sigma \in \Sigma$ .
- for**  $i = 1$  **to**  $m$
- Compute the rank  $R_i$  truncated SVD  $\hat{\mathcal{H}}_{(1),i} \simeq \mathbf{U}_i \mathbf{D}_i \mathbf{V}_i^\top$ .
- Let  $A_i = (\mathbf{U}_i^\top \mathbf{U} \alpha, \{\mathbf{U}_i^\top \mathbf{U} \mathbf{A}^\sigma \mathbf{U}_i^\top \mathbf{U}\}_{\sigma \in \Sigma}, \mathbf{U}_i^\top \mathbf{U} \Omega_{:,i})$

- Additional step** to the spectral learning algorithm (lines 3-5):
  - vv-WFA  $A = (\alpha, \{\mathbf{A}^\sigma\}_{\sigma \in \Sigma}, \Omega)$  is minimal  $\nRightarrow$  WFA  $A_i = (\alpha, \{\mathbf{A}^\sigma\}_{\sigma \in \Sigma}, \Omega_{:,i})$  is minimal.
  - Need to **project down each  $A_i$**  to its true dimension.

## Theoretical Insight

**Theorem** Let  $\mathbf{M} \in \mathbb{R}^{d_1 \times d_2}$  of rank  $R$ ,  $\hat{\mathbf{M}} = \mathbf{M} + \mathbf{E}$ ,  $\Pi_U, \Pi_V \in \mathbb{R}^{d_1 \times d_1}$  matrices of orthogonal projections on the top  $R$  left sing. vectors of  $\mathbf{M}$  and  $\hat{\mathbf{M}}$ . Then, for any  $\delta > 0$ , with probability  $\geq 1 - \delta$ ,

$$\|\Pi_U - \Pi_V\|_F \leq 4 \left( \sqrt{\frac{(d_1 - R)R + 2 \log(1/\delta)}{d_1 d_2}} \frac{\|\mathbf{E}\|_F}{s_R(\mathbf{M})} + \frac{\|\mathbf{E}\|_F^2}{s_R(\mathbf{M})^2} \right). \quad (1)$$

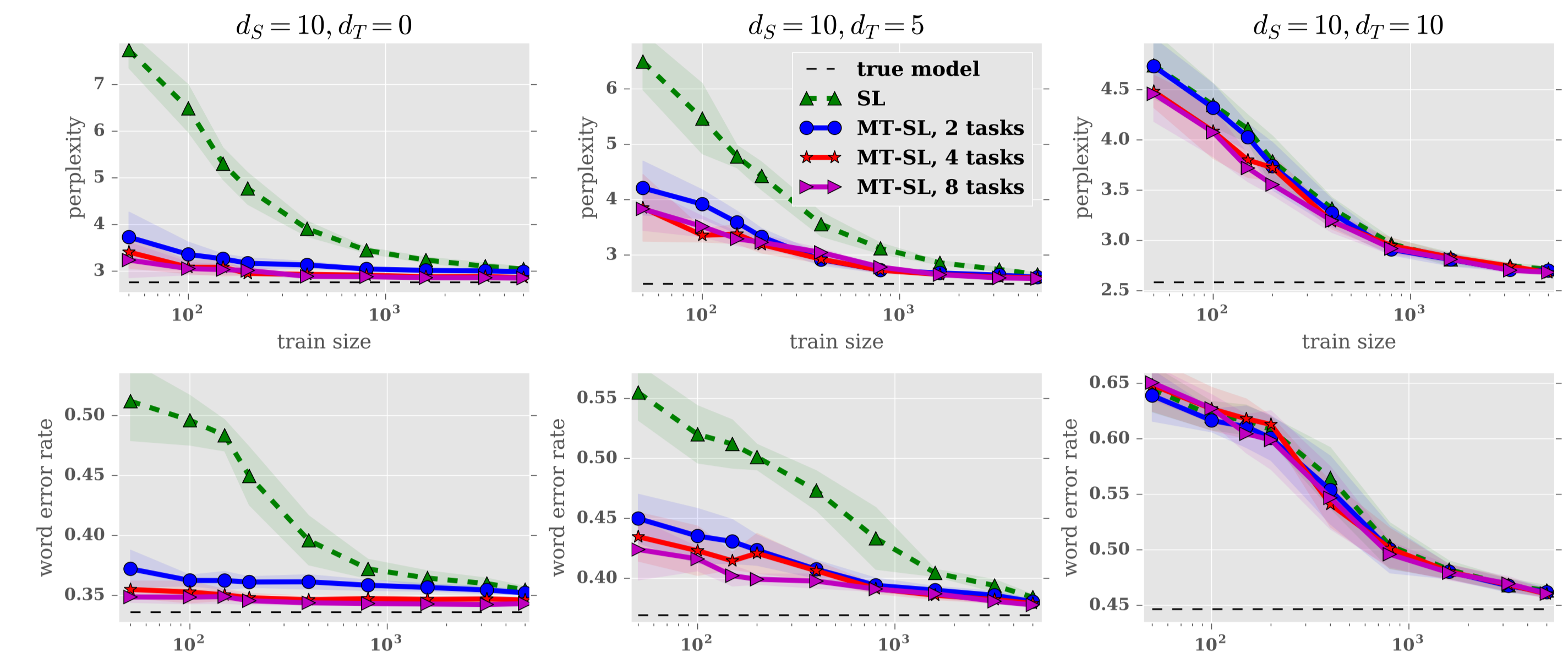
- Consider  $m$  tasks  $f_1, \dots, f_m$  with empirical Hankel matrices  $\hat{\mathbf{H}}_1, \dots, \hat{\mathbf{H}}_m \in \mathbb{R}^{P \times S}$ , then

$$\hat{\mathcal{H}}_{(1)} = [\hat{\mathbf{H}}_1 \hat{\mathbf{H}}_2 \dots \hat{\mathbf{H}}_m] \in \mathbb{R}^{P \times mS}.$$

- If the tasks are maximally related (i.e.  $R = \text{rank}(\vec{f}) = \text{rank}(f_1) = \dots = \text{rank}(f_m)$ ) then as the number of tasks grows, the first term in Eq. (1) tends to 0 and **the estimation error of the singular subspace decays quadratically instead of linearly**.

## Experiments on Synthetic Data

- Randomly generated stochastic WFAs following the PAutoMaC competition process [6].
- Related WFAs**: joint feature space of dimension  $d_S = 10$  and task specific space of dimension  $d_T$  (i.e.  $\text{rank}(f_i) = d_S + d_T$  and  $\text{rank}(\vec{f}) = \text{rank}([f_1, \dots, f_m]) = d_S + md_T$ ).
- Training sample drawn from target task  $f_1$  and training samples of size 5,000 for tasks  $f_2, \dots, f_m$ .



## Experiments on Real Data

- Universal Dependencies [5]: sentences from 33 languages labeled with 17 PoS tags.
- Samples drawn from **33 distributions** over strings on an **alphabet of size 17**.
- For each language, (80%, 10%, 10%)-split between training, validation and test sets.
- Two ways of selecting related tasks:
  - use all other languages
  - select the 4 closest languages w.r.t. the distance between the (top-50) left singular subspaces of the Hankel matrices.

Training size	100	500	1000	5000	all available data
Related tasks: all other languages					
Perplexity	7.0744 (±7.76)	3.6666 (±5.22)	3.2879 (±5.17)	3.4187 (±5.57)	3.1574 (±5.48)
WER	1.4919 (±2.37)	1.3786 (±2.94)	1.2281 (±2.62)	1.4964 (±2.70)	1.4932 (±2.77)
Related tasks: 4 closest languages					
Perplexity	6.0069 (±6.76)	4.3670 (±5.83)	4.4049 (±5.50)	2.9689 (±5.87)	2.8229 (±5.90)
WER	2.0883 (±3.26)	1.5175 (±2.87)	1.2961 (±2.57)	1.3080 (±2.55)	1.2160 (±2.31)

Table: Average relative improvement over all languages (in %) of MT-SL vs. SL on the UNIDEP dataset (e.g. for perplexity we report  $100 \cdot (\rho_{\text{SL}} - \rho_{\text{MT-SL}}) / \rho_{\text{SL}}$ ).

Target task	4 closest tasks w.r.t. subspace distance (closest first)			
Basque	Finnish	Polish	Czech	Indonesian
Croatian	Estonian	Slovenian	Czech	Finnish
French	Italian	Spanish	German	English
Hungarian	Danish	Ancient Greek	German	Portuguese
Gothic	Old Church Slavonic	Latin	Ancient Greek	Finnish
Italian	English	French	Spanish	Dutch
Japanese	Hindi	Persian	Arabic	Tamil
Latin	Old Church Slavonic	Ancient Greek	Gothic	Finnish
Swedish	Danish	Norwegian	Finnish	Estonian

Table: Some related tasks used in the UNIDEP experiment.

- Cherry picked example: on the Basque task with a training set of size 500, the **WER was reduced from ~76% for SL to ~70% using all other languages as related tasks, and to ~65% using the 4 closest tasks** (Finnish, Polish, Czech and Indonesian).

## References

- Borja Balle, Xavier Carreras, Franco M Luque, and Ariadna Quattoni. Spectral learning of weighted automata. *Machine learning*, 96(1-2):33–63, 2014.
- T Tony Cai and Anru Zhang. Rate-optimal perturbation bounds for singular subspaces with applications to high-dimensional statistics. *arXiv preprint arXiv:1605.00353*, 2016.
- Jack W. Carlyle and Azaria Paz. Realizations by stochastic finite automata. *Journal of Computer and System Sciences*, 5(1):26–40, 1971.
- Michel Fliess. Matrices de hankel. *Journal de Mathématiques Pures et Appliquées*, 53(9):197–222, 1974.
- Joakim Nivre, Zeljko Agić, Lars Ahrenberg, et al. Universal dependencies 1.4, 2016. URL <http://hdl.handle.net/11234/1-1827>. LINDAT/CLARIN digital library at the Institute of Formal and Applied Linguistics, Charles University.