

Differentiable Dynamic Programming for Structured Prediction and Attention



Arthur Mensch⁽¹⁾

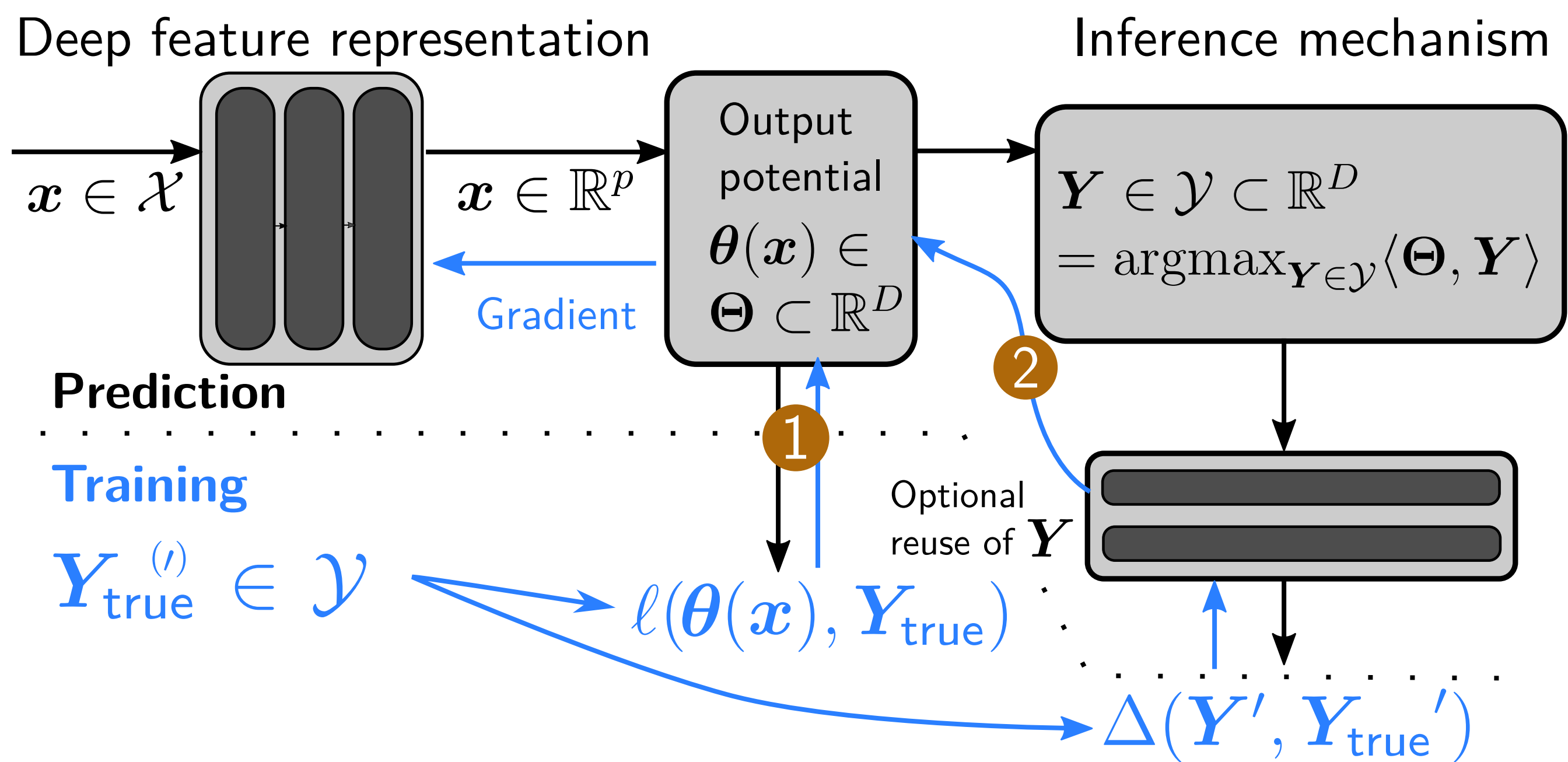
Mathieu Blondel⁽²⁾

⁽¹⁾Parietal team, Inria, Saclay, France

⁽²⁾NTT Communication Science Laboratories, Kyoto, Japan



1 – Context – End-to-end deep structured prediction



- $x \in \mathcal{X}$: text sequence, audio time-series
- $Y \in \mathcal{Y}$: tag sequence, alignment between time-series, **attention map**
- Losses typically involve LP solution — baseline structured perceptron

$$l(\theta(x), Y_{\text{true}}) = \max_{Y \in \mathcal{Y}} \langle Y, \theta \rangle - \langle Y_{\text{true}}, \theta \rangle$$

- Easiest case: $\theta(x) = Wx + b$, $\mathcal{Y} = [1, k]$
- Often in machine learning: max and argmax from dynamic programming

Dynamic programming

End-to-end training

$$(\text{arg})\max_{Y \in \mathcal{Y}} \langle \theta(x), Y \rangle$$

breaking \mathcal{Y} into smaller sets

- Linear vs exp complexity
- Inference on tree CRF / DTW

- 1 Loss functions
⇒ Backpropagate through LP value (max)
- 2 Differentiable inference mechanism
⇒ Backpropagate through LP solution (argmax)

2 – Contrib – Smoothing DP for end-to-end training

- **Dynamic programming value/solution become differentiable layer**
- End-to-end training for many supervised prediction problems
- Easy to implement custom backward module, with probabilistic interpreta^o
- Optional **sparsity** in $Y \Rightarrow$ small set of best guesses / sparse attention

3 – Generic DP on a directed acyclic graph

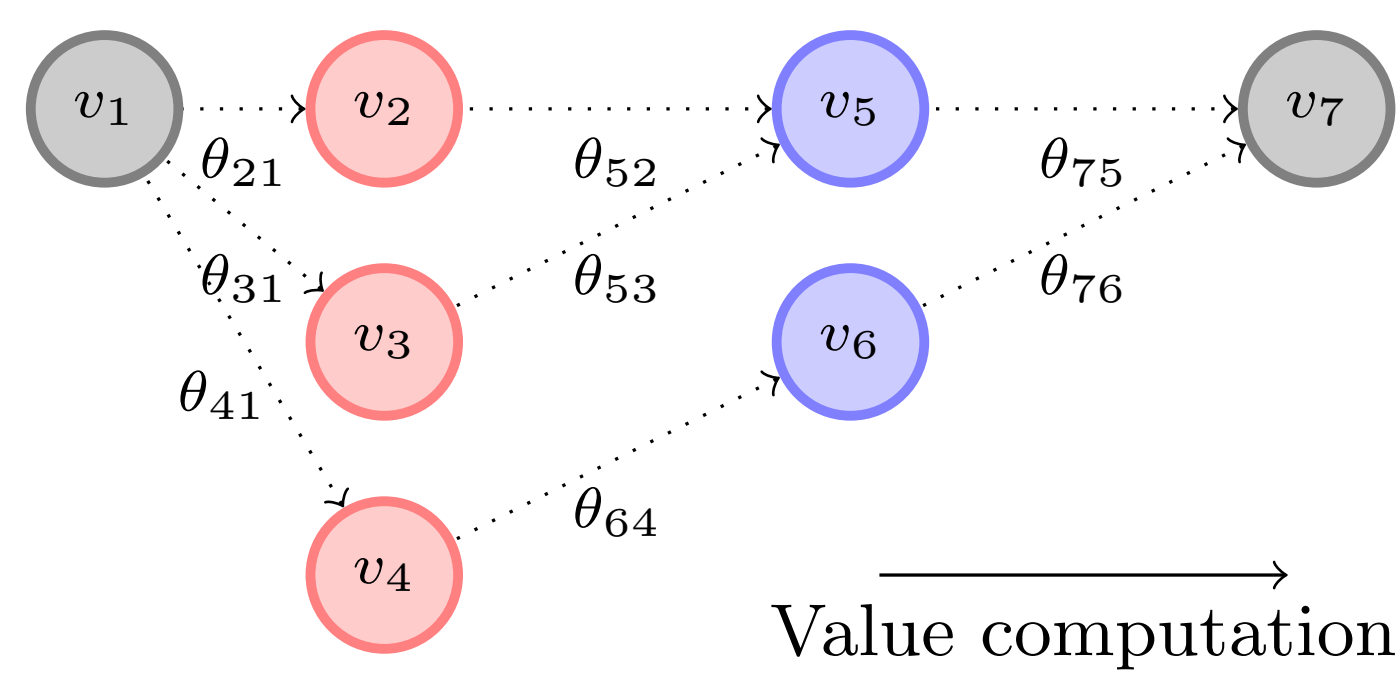
- Find shortest path on a graph G . Edge (i, j) has weight θ_{ij} — j parent, i child
- Path $Y \in \mathcal{Y} \subset \{0, 1\}^{N \times N}$: $y_{i,j} = 1$ iff (i, j) is taken, $\theta \in \mathbb{R}^{N \times N}$ adj. matrix

Bellman equation

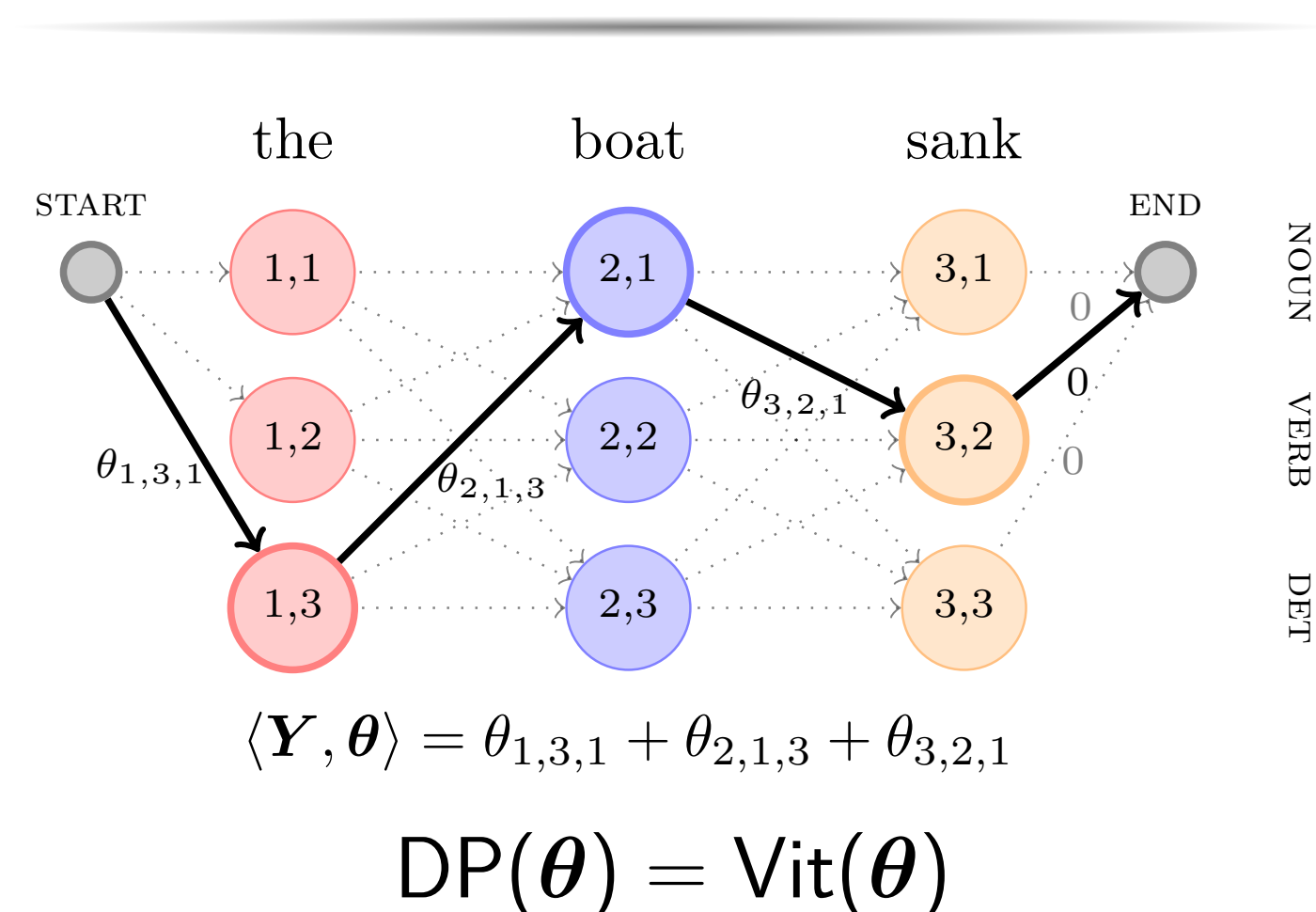
$$v_i(\theta) = \max_{j \in P_i} \theta_{i,j} + v_j(\theta), DP(\theta) \triangleq v_n(\theta)$$

Highest score among all paths:

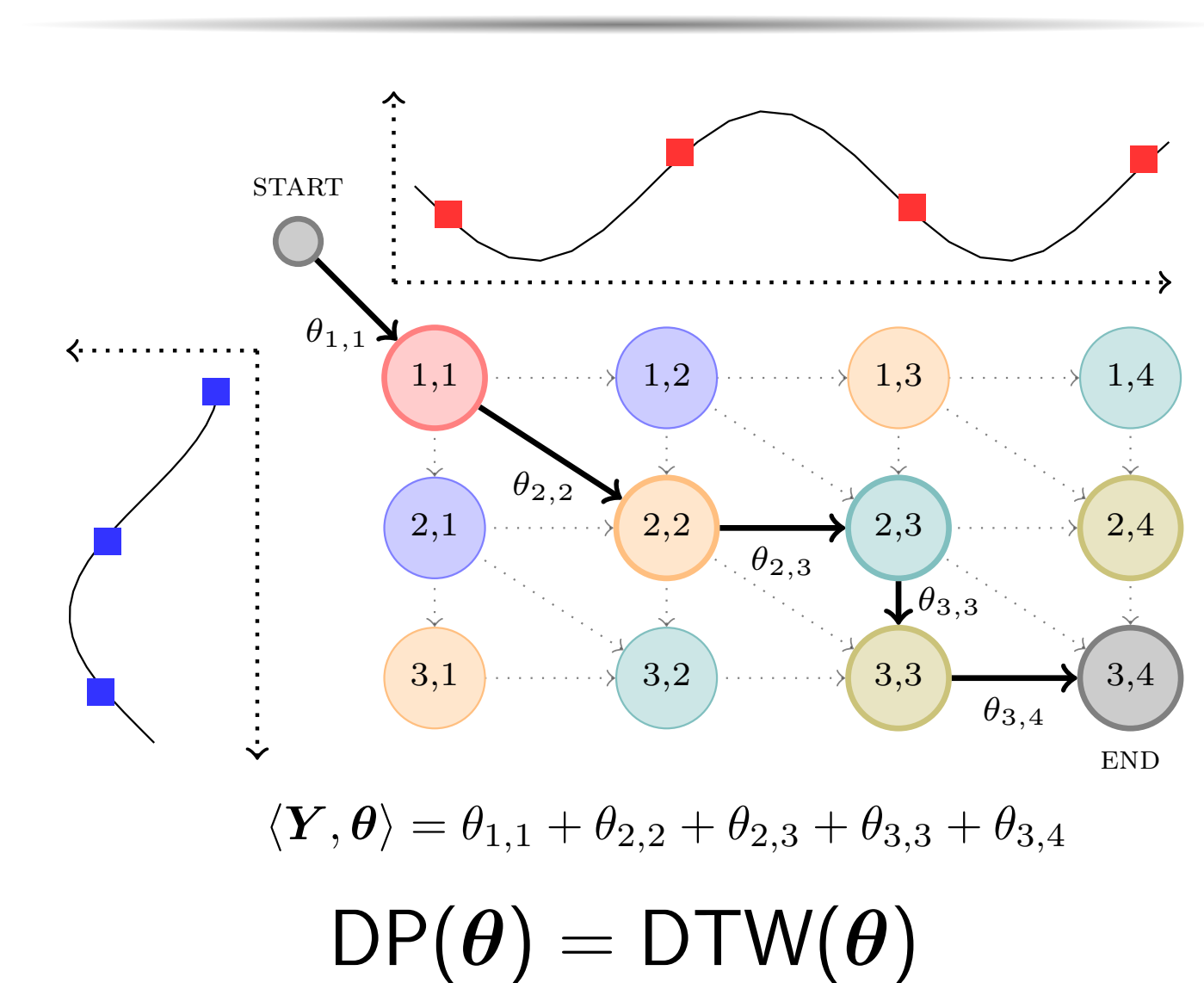
$$LP(\theta) = \max_{Y \in \mathcal{Y}} \langle Y, \theta \rangle = DP(\theta)$$



Viterbi decoder



Dynamic time warping



Backpropagation

Danskin theorem

$$\partial DP(\theta) = \text{conv}(\text{argmax}_{Y \in \mathcal{Y}} \langle Y, \theta \rangle)$$

- Not defined everywhere, not smooth, piecewise constant
- Need regularization

4 – DP smoothing

- $\Omega : \mathbb{R} \rightarrow \mathbb{R}$ str-convex f (Nesterov, 2005).
 $\max_{\Omega}(x) = \max_{y \in \Delta^d} \langle x, y \rangle - \sum_{i=1}^d \Omega(y_i)$
- **Bottom-up Bellman smoothing**
 $DP(\theta) \rightarrow DP_{\Omega}(\theta)$, $\partial DP(\theta) \rightarrow \nabla DP_{\Omega}(\theta)$
- Vs generally intractable $LP_{\Omega}(\theta)$
- Argmax backprop: $\nabla^2 DP_{\Omega}(\theta)$

Regularizations have different properties

- **Negentropy** $-H$: $\Omega(p) = p \log(p)$
- $\nabla \max_{\Omega}(x) = \text{softmax}(x)$
- ∇DP_{Ω} expected Gibbs distribution
- ℓ_2^2 : $\Omega(p) = p^2$ (Martins et al., 2016)
- $\nabla \max_{\Omega}(x) = \text{sparsemax}(x)$
- Projection on Δ^d : ∇DP_{Ω} **sparse**

5 – Theoretical properties

DP value layer

- DP_{Ω} is convex, LP – DP_{Ω} bounded
- $DP_{\Omega}(\theta) = LP_{\Omega}(\theta)$ if and only if $\Omega = -\gamma H$, where $\gamma \geq 0$

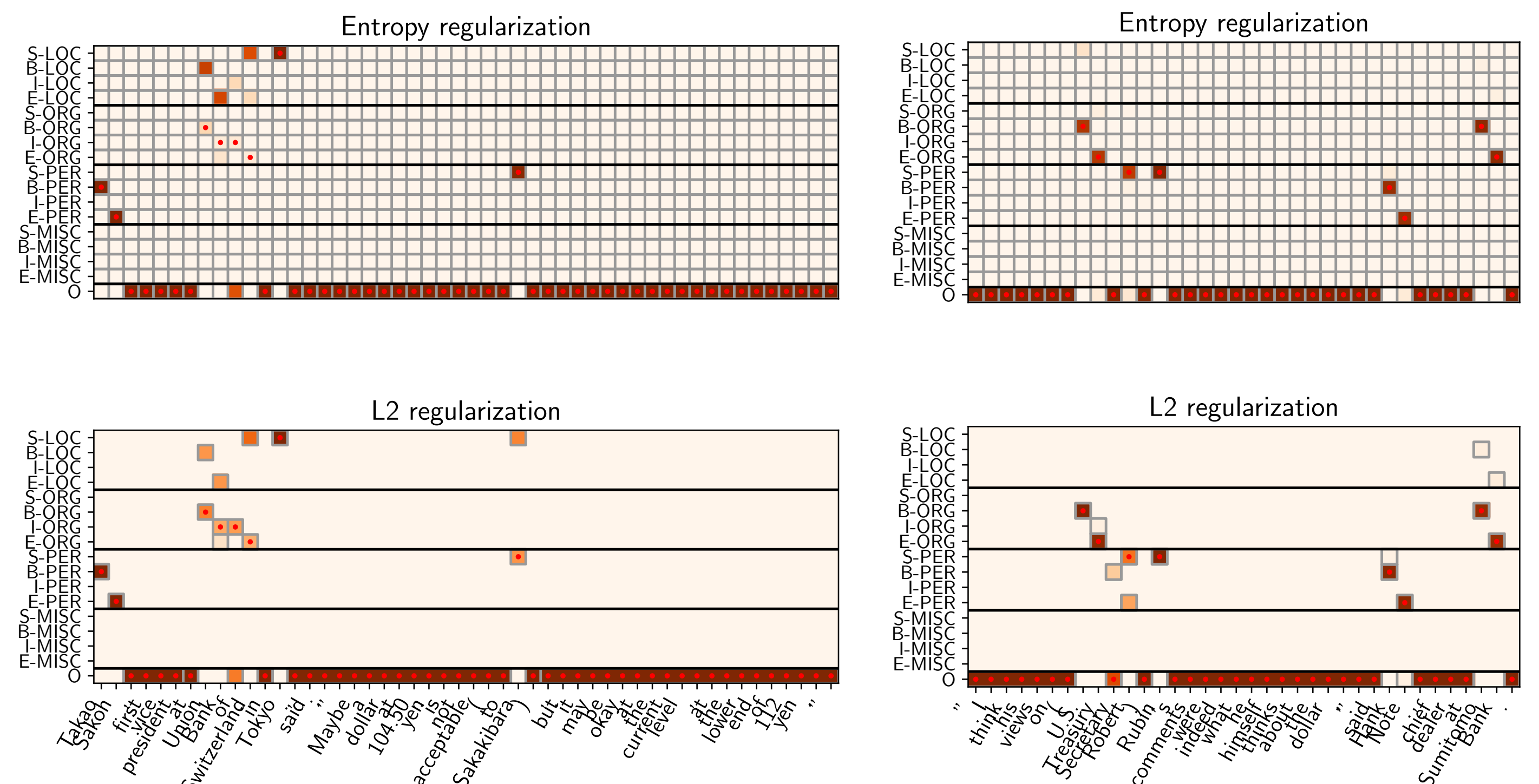
DP gradient layer

- Path distribution $p_{\theta, \Omega}$ (over \mathcal{Y}) s.t.
 $\nabla DP_{\Omega}(\theta) = \mathbb{E}_{p_{\theta, \Omega}}[Y] \in \text{conv}(\mathcal{Y})$
= Graph random walk, easy to sample

6 – Experiments – Named entity recognition (CoNLL'03)

- Word embedding + word LSTM
- Linear chain CRF-like potential
 $\theta(X)_{t,i,j} = w_i^T x_t + b_i + t_{i,j}$
- Sparse prediction $\nabla DP_{\ell_2^2}(\theta)$

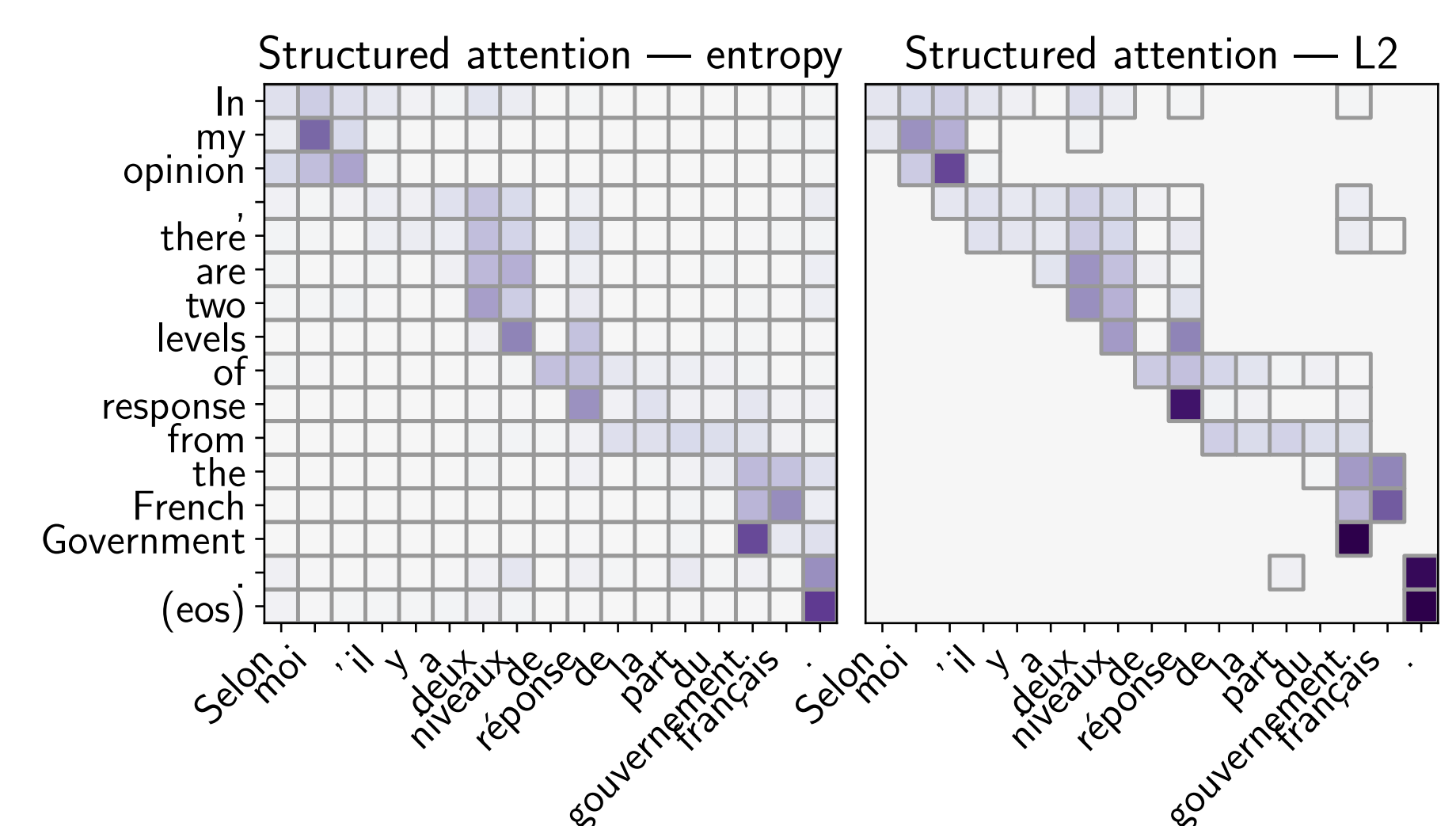
Ω	Loss	English	Spanish	German	Dutch
Negent.	Surrogate	90.80	86.68	77.35	87.56
	Relaxed	90.47	86.20	77.56	87.37
ℓ_2^2	Surrogate	90.86	85.51	76.01	86.58
	Relaxed	89.49	84.07	76.91	85.90
0	Struct. perceptron	86.52	81.48	68.81	80.49
Lample et al., 2016		90.96	85.75	78.76	81.74



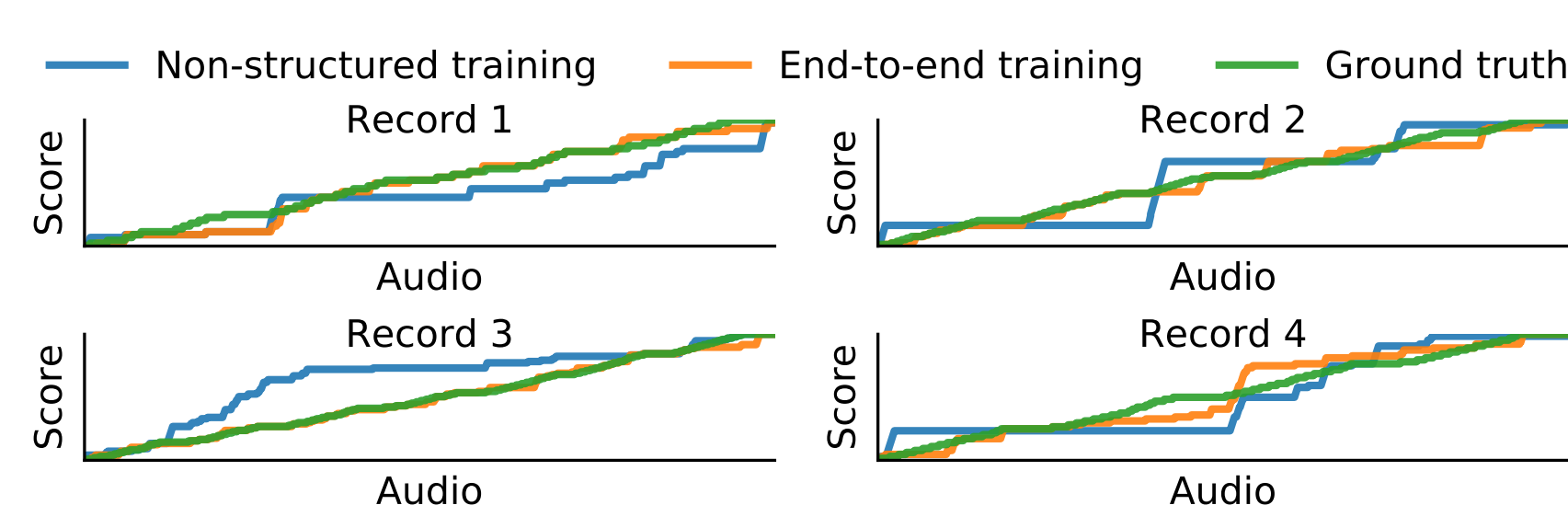
Exp – Structured and sparse attention (WMT14 fr-en)

- Compute attention c by marginalizing a graphical model (Vit_{Ω}), with sparse marginal computation $\Omega = \ell_2^2$. vs simple softmax in original version
- Regularization is needed to backprop through argmax
- Vit_{Ω} double-backward using Pearlmutter (1994) AD

Attention	WMT14 IM BLEU
fr → en Softmax	27.96
fr → en Entropy reg.	27.96
ℓ_2^2 reg.	27.21
en → fr Softmax	28.08
en → fr Entropy reg.	27.98
ℓ_2^2 reg.	27.28



Exp – Audio to score alignment



7 – Code + further

- PyTorch + AllenNLP implem:
github.com/arthurmensch/didyprog
- More on algo + theory in paper
- See also SparseMAP (Nicolae et al.) for a different approach at ICML'18

Lample, Guillaume et al. (2016). "Neural Architectures for Named Entity Recognition". In: *Proc. of NAACL*, pp. 260–270.
 LeCun, Yann et al. (2006). "A tutorial on energy-based learning". In: *Predicting structured data 1.0*.
 Martins, André F.T. and Ramón Fernández Astudillo (2016). "From softmax to sparsemax: A sparse model of attention and multi-label classification". In: *Proc. of ICML*, pp. 1614–1623.
 Mensch, Arthur and Mathieu Blondel (2018). "Differentiable Dynamic Programming for Structured Prediction and Attention". In: *Proceedings of the International Conference on Machine Learning*.
 Nesterov, Yurii (2005). "Smooth minimization of non-smooth functions". In: *Mathematical Programming* 103.1, pp. 127–152.
 Nicolae, Vlad et al. (2018). "SparseMAP: Differentiable sparse structured inference". In: *Proceedings of the International Conference on Machine Learning*.
 Pearlmutter, Barak A. (1994). "Fast exact multiplication by the Hessian". In: *Neural computation* 6.1, pp. 147–160.