

1 We kindly thank the reviewers for their feedback. All reviewers expressed their interest in more experimental evaluation.
2 We report an additional experiment on sequence prediction and then address each reviewer’s questions individually.

3 **Sequence Prediction.** We considered the *TIMIT Acoustic-Phonetic Continuous Speech Corpus* (<https://catalog.ldc.upenn.edu/LDC93S1>), containing recordings of ~ 6300 recorded sentences (630 English speakers, each reading
4 ~ 10 sentences). We preprocessed the dataset following a standard practice in the speech community, taking 10ms
5 frames (dropping the glottal stop ‘q’ labeled frames in the core test set) and mapping them as 40-dimensional vectors [
6 D. Povey et al., “*The kaldı speech recognition toolkit*”, 2011]. The speech recognition problem is stated in terms of
7 sequence prediction, i.e. each sentence, represented as a list of 40-dimensional vectors to be mapped in a same-length
8 list of phonemes (131 phonemes per each 10ms frame). We considered:
9

- 10 1. A baseline performing independent classification of sequence elements using Kernel SVM with Gaussian Kernel.
11 The regularization parameter C and the bandwidth of the kernel have been selected via hold-out cross validation in
12 the logarithmic range $[10^{-3}, 10^3]$ and $[10^{-9}, 10^3]$.
- 13 2. Structural SVM for sequence tagging based on Hidden Markov Models (SVM-HMM) and exact Viterbi algorithm
14 for inference (see below). Implementation from https://www.cs.cornell.edu/people/tj/svm_light/svm_hmm.html using Gaussian kernel. The regularization parameter C and the bandwidth of the kernel have been
15 selected via hold-out cross validation in the logarithmic range $[10^{-3}, 10^3]$ and $[10^{-9}, 10^3]$.
- 16 3. The proposed *localized structured prediction algorithm*. In particular we defined parts on input as subsequences of
17 11 frames (110ms) corresponding to vectors of 440 elements, while the corresponding parts of the output are the
18 couple of phonemes associated to the central and the next frame of the input.
19

20 *Inference.* For both SVM-HMM and our method, the inference is solved with the classic Viterbi algorithm over a table
21 of dimension $L|x|$ where L is the number of labels per element ($L = 131$ phonemes) and $|x|$ is the length of the input
22 sequence ($|x| \sim 200$ on average). When output parts are couples of phonemes this entails $O(|x|L^2)$ computations.
23

24 *Results.* The Table reports the test performance (Hamming loss)
25 of the three methods over 5 trials (dataset is reshuffled and
26 randomly split in 3300 sentences for training and the rest for
27 test). Leveraging the parts structure appears beneficial, with
28 our method consistently outperforming both the baseline and
29 the traditional structured prediction competitor.

Sequence prediction on TIMIT	Hamming
Independent SVM	0.31 ± 0.015
Struct SVM	0.29 ± 0.031
Localized Struct. Pred	0.26 ± 0.012

30 **R1** * *Long Range Dependencies.* Our assumption does not limit long-range dependencies in the data as they do not
31 imply that two parts to have identical (or very similar) appearance with respect to the similarity metric used. However,
32 it is also true that while such dependencies could be leveraged to improve performance, our approach is not designed to
33 capture them at training time (but only during inference). As future work we will explore the question of reformulating
34 the part-based regression as a multi-task learning (MTL) problem, with each task corresponding to a part. Then,
35 by leveraging ideas from the MTL literature we could impose (when known a-priori) or learn (when unknown) the
36 multi-scale relations/dependencies between tasks/parts during the training phase. Algorithmically this extension would
37 require a small modification of the current approach. We will add a comment on this when discussing future directions.

38 * *Inference.* The inference step in our setting is formulated as an optimization problem over the output space Z as most
39 traditional structured prediction approaches (see Remark 2 in our paper). Hence, our inference step is as difficult as
40 previous methods, from a computational viewpoint (while also providing strong theoretical guarantees on the prediction).
41 For instance, in the experiments on TIMIT in this document, both our method and SVM-HMM use the viterbi algorithm
42 for inference. Algorithm 3 offers an additional benefit of our formulation when first order optimization on Z is possible.

43 **R2** * *Example 2.* We agree with R2 to improve exposition of Example 2. In particular we will clarify the difference
44 between the loss function used to evaluate the error and the score function at inference time. While these are two
45 separate concepts in the context of CRFs, they essentially coincide in our Localized Structured Prediction framework
46 (score is a linear combination of losses). The conclusion of Example 2 highlights the connection between CRFs and
47 our framework by observing that the score functions have essentially the same structure. Hence they lead to the same
48 inference problem. [Sutton, C and McCallum, A. “*An introduction to conditional random fields.*” 2012].

49 * *Notation* We agree with R2: will use $\Delta(z, y)$ in the main text and $\Delta(z, y|x)$ only in the appendix. This notation was
50 originally introduced to account for settings where the parts might depend on the input (e.g. images with different sizes).

51 **R3** * *Attention Models.* Loosely speaking, an attention model could be formulated within our framework by considering
52 a parametrization of the possible part structures: in this sense, the attention would consist in the process of selecting
53 iteratively the parts that are more relevant to the task and adapt them depending on individual inputs. This model
54 however would require to perform an optimization over the parts parametrization, making both learning and inference
55 significantly more challenging. Investigating this question could lead to interesting future work, we thank the reviewer
56 and we will add a comment to the paper. * *Code.* Our code is in python (+ pytorch) and will be made available in Github.