

# MUST-CNN: A Multilayer Shift-and-Stitch Deep Convolutional Architecture for Sequence-based Protein Structure Prediction

**Zeming Lin**

Department of Computer Science  
University of Virginia  
Charlottesville, VA 22904  
zl4ry@virginia.edu

**Jack Lanchantin**

Department of Computer Science  
University of Virginia  
Charlottesville, VA 22904  
jjl5sw@virginia.edu

**Yanjun Qi**

Assistant Professor  
Department of Computer Science  
University of Virginia  
Charlottesville, VA 22904  
yq2h@virginia.edu

## Abstract

Predicting protein properties such as solvent accessibility and secondary structure from its primary amino acid sequence is an important task in bioinformatics. Recently, a few deep learning models have surpassed the traditional window based multilayer perceptron. Taking inspiration from the image classification domain we propose a deep convolutional neural network architecture, MUST-CNN, to predict protein properties. This architecture uses a novel multilayer shift-and-stitch (MUST) technique to generate fully dense per-position predictions on protein sequences. Our model is significantly simpler than the state-of-the-art, yet achieves better results. By combining MUST and the efficient convolution operation, we can consider far more parameters while retaining very fast prediction speeds. We beat the state-of-the-art performance on two large protein property prediction datasets.

## Introduction

Proteins are vital to the function of living beings. It is easy to determine the sequence of a protein, yet it is difficult to determine other properties, such as secondary structure and solvent accessibility. These properties are hypothesized to be almost uniquely determined by primary structure, but it is still computationally difficult to determine them on a large scale.

Previous state-of-the-art methods for protein secondary structure prediction use multilayer perceptron (MLP) networks (Qi et al. 2012; Drozdetskiy et al. 2015). In order to predict a per-position label for each amino acid in the input protein sequence, MLP networks must use a “windowing” approach where a single label is predicted by feeding the target amino acid and its surrounding amino acids through the network. This is then repeated for each amino acid in the sequence. These architectures generally has two major drawbacks due to the windowing approach: (1) they take a long time to train, on the order of days or weeks and (2) they have smaller window sizes, and thus cannot make longer range connections. For instance, the PSIPRED algorithm handles a window size of only 15 amino acids (Jones 1999).

To overcome the windowing issue, we propose to use a convolutional neural network (CNN) which can label the properties of each amino acid in the entire target sequence all at once. CNNs have been used successfully in computer vision (Pinheiro and Collobert 2013; Szegedy et al. 2014) and natural language processing (Kim 2014; Collobert and Weston 2008). In addition to parameter sharing and pooling, which reduce computation, CNNs are also highly parallelizable. Thus, CNNs can achieve a much greater speedup compared to a windowed MLP approach. The issue when trying to label each position in an input sequence with a CNN is that pooling leads to a decreased output resolution. To handle this issue, we propose a new multilayer shift-and-stitch method which allows us to efficiently label each target input at full resolution in a computationally efficient manner.

We show that a Multilayer Shift-and-stitch CNN (MUST-CNN) trained end-to-end and per-position on protein property prediction beats the state-of-the-art without other machinery. To our knowledge, we are the first to train convolutional networks end-to-end for per-position protein property prediction. Both learning and inference are performed on entire arbitrarily sized sequences. Feedforward computation and backpropagation is made possible by our novel application of the shift-and-stitch technique on the entire sequence.

In summary we make the following contributions:

1. Beat the state-of-the-art performance on two large datasets of protein property prediction tasks.
2. Propose a multilayer shift-and-stitch technique for deep CNNs, which significantly speeds up training and test time and increases the size of the model we can train.
3. Propose a generic end-to-end system for per-position labeling on the sequence level. That is, for a sequence  $\{a_k\}_{k=1}^n$ , we can generate labels  $\{y_k\}_{k=1}^n$  for each  $a_k$ .

## Related Works

Two of the most used algorithms in bioinformatics literature for protein property prediction are PSIPRED (Jones 1999) and Jpred (Drozdetskiy et al. 2015). PSIPRED 3.2, which uses a two layer MLP approach, claims a 3-class per-position accuracy ( $Q_3$ ) score of 81.6%. The Jpred algorithm uses a very similar structure of a two layer MLP network. However, Jpred considers more features and uses

a jury based approach with multiple models (Cuff and Barton 2000). Jpred claims an 81.5%  $Q_3$  score on secondary structure prediction, and also predicts relative solvent accessibility. (Qi et al. 2012) uses a deep MLP architecture with multitask learning and achieves 81.7%  $Q_3$ . (Zhou and Troyanskaya 2014) created a generative stochastic network to predict secondary structure from the same data we used, for a  $Q_8$  of 66.4%. Unlike  $Q_3$ , the  $Q_8$  accuracy tries to distinguish between more classes.

The state-of-the-art protein sequence classification system is SSpro, which obtains 91.74%  $Q_3$  and 85.88%  $Q_8$  on a different unfiltered PDB dataset (Magnan and Baldi 2014). However, this system exploits additional information via sequence similarity, and their reported accuracies were only 80% without this module. Our work would complement their machine learning module and likely result in even better accuracies.

Recently, work has also been done on the model side, particularly in natural language processing and image recognition tasks. (Collobert et al. 2011) created a similar algorithm in the natural language processing domain, where they labeled word properties, such as part of speech or category of a named entity, on text data. If we consider each protein chain to be a sentence and each amino acid to be a word, the techniques transfer easily. (Collobert et al. 2011) used both a windowed approach and a sentence level approach with a convolutional network, though their network was shallow and only outputted predictions for one position at a time. Long-short term memory networks have been used very successfully in sequence learning, machine translation (Sutskever, Vinyals, and Le 2014; Bahdanau, Cho, and Bengio 2014) and language modeling (Zaremba, Sutskever, and Vinyals 2014). We note that machine translation is a much more general sequence to sequence task where the input and output sizes are not matched. Language modeling tries to guess future words based on past words, while protein sequences has no innate direction.

In the image domain, (Szegedy et al. 2014) has beaten the state-of-the-art on image classification by a large percentage through using a deep multilayer convolutional network in the ImageNet Large-Scale Visual Recognition Challenge. Scene labeling is the task of labeling each pixel of an image with one of several classes, a 2D analogue of protein property prediction. (Pinheiro and Collobert 2013) uses a recurrent neural network to obtain state-of-the-art results on scene labeling without any feature engineering. (Long, Shelhamer, and Darrell 2014) designs fully convolutional networks for dense pixel prediction by running several convolutional networks on different scales. (Sermanet et al. 2013) increases the resolution of a bounding box based image classifier by introducing the shift-and-stitch technique, which we use on sequences instead of images and on the entire model instead of only on the last layer.

## Method: MUST-CNN

### Convolutional Neural Networks (CNN)

Convolutional networks were popularized for the task of handwriting recognition of 2D images (Lecun et al. 1998).

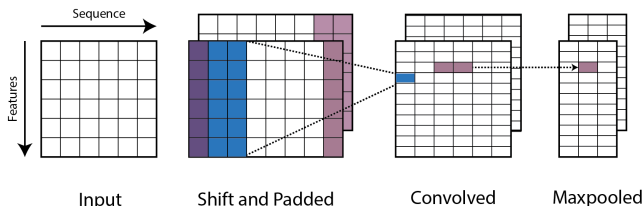


Figure 1: A diagram for one layer of the convolutional network. We shift and pad the input accordingly to be able to label every position in our neural network. How the output is recombined into a fully dense per-position sequence label is given in Figure 3.

In a similar way, we use a 1D convolution for the protein sequence labeling problem. A convolution on sequential data tensor  $\mathbf{X}$  of size  $T \times n_{in}$  with length  $T$ , kernel size  $k$  and input hidden layer size  $n_{in}$  has output  $\mathbf{Y}$  of size  $T \times n_{out}$ :

$$\mathbf{Y}_{t,i} = \sigma(\mathbf{B}_i + \sum_{j=1}^{n_{in}} \sum_{z=1}^k \mathbf{W}_{i,j,k} \mathbf{X}_{t+z-1,j})$$

where  $\mathbf{W}$  and  $\mathbf{B}$  are the trainable parameters of the convolution kernel, and  $\sigma$  is the nonlinearity. We try three different nonlinearity functions in our experiments: the hyperbolic tangent, rectified linear units (ReLU), and piecewise rectified linear units (PReLU). The hyperbolic tangent is historically the most used in neural networks, since it has nice computational properties that make optimization easy. Both ReLU and PReLU have been shown to work very well on deep convolutional networks for object recognition. ReLU was shown to perform better than tanh on the same tasks, and enforces small amounts of sparsity in neural networks (Glorot, Bordes, and Bengio 2011). By making the activations trainable and piecewise, PReLU has shown to match the state of the art on ILSVRC while converging in only 7% of the time (He et al. 2015).

The ReLU nonlinearity is defined as

$$\text{relu}(x) = \max(0, x)$$

and the PReLU nonlinearity is defined as

$$\text{prelu}(x) = \begin{cases} \alpha x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$

with a trainable parameter  $\alpha$ .

After the convolution and nonlinearity, we use a pooling layer. The only pooling strategy tested was maxpooling, which has shown to perform much better than subsampling as a pooling scheme (Scherer, Müller, and Behnke 2010) and has generally been the preferred pooling strategy for large scale computer vision tasks. Maxpooling on a sequence  $\mathbf{Y}$  of size  $T \times n$  with a pooling size of  $m$  results in output  $\mathbf{Z}$  where

$$\mathbf{Z}_{t,i} = \max_{j=1}^m \mathbf{Y}_{m(t-1)+j,i}$$

Finally, the outputs are passed through a dropout layer. The dropout layer is a randomized mask of the outputs, equivalent to randomly zeroing out the inputs to the next layer during training time with probability  $d$  (Srivastava et al. 2014). During testing, the dropout layer is removed and all weights are used. This acts as a regularizer for the neural network and prevents overfitting, though the best values for  $d$  must be discovered experimentally.

One layer of the convolutional network is depicted in Figure 1. In our model design, we apply the CNN module multiple times for a deep multilayer framework.

### Multilayer Shift-and-Stitch (MUST)

Pooling is a dimension reduction operation which takes several nearby values and combines them into a single value – maxpooling uses the max function to do this. Maxpooling is important because as nearby values are merged into one, the classifier is encouraged to learn translation invariance. However, after a single application of maxpooling with a pool size of  $m$  on input sequence  $\mathbf{X}$  of length  $T$ , the resulting maxpool output has sequence length  $\frac{T}{m}$ .

Since the dimensionality of the sequence has been divided by a factor of  $m$ , it is no longer possible to label every position of the original sequence. A technique to increase the resolution in convolutional networks was given in (Sermanet et al. 2013), called “shift-and-stitch”. Their implementation uses the technique in a two dimensional setting to increase the resolution of pixel labels in the last layer of a convolutional network, for up to a  $4\times$  increase in resolution. We observe that the limiting factor on applying this to an entire image is the massive slowdown in computation, since each pooling layer in a two-dimensional case requires the network to stitch together 4 different outputs and 3 pooling layers require 64 different stitched outputs.

However, in the sequential case, we need to stitch together significantly fewer sequences. Using 3 pooling layers with pooling size 2 will only require 8 different stitches, making computation tractable. Therefore, we propose to apply shift-and-stitch to every layer of our deep CNN which generates dense per-position predictions for the entire sequence. This process is described in Figure 3. This will allow us to take advantage of the computational speeds provided by the convolution module, making it feasible to try a much larger model.

Due to the kernel size, a convolution with kernel size  $k$  removes the  $\lfloor \frac{k}{2} \rfloor$  edge values on each end of the sequence. Thus, we pad the input with a total of  $\lfloor \frac{k}{2} \rfloor - 1$  zeros at each end, colored as red in Figures 1 and 3. Because a maxpooling operation with pooling size  $m$  labels every  $m$  values in the input, we duplicate the input  $m$  times and pad the  $i$ -th input such that the first convolution window is centered on the first amino acid. We observe that we can then join the  $m$  duplicated inputs along the batch dimension and pass it into the convolution module and take advantage of the batch computation ability offered by standard linear algebra packages to train our system even faster. After pooling, the output is a zipped version of the original input along the batch dimension. We simply “stitch” together the output in full resolution for the final result.

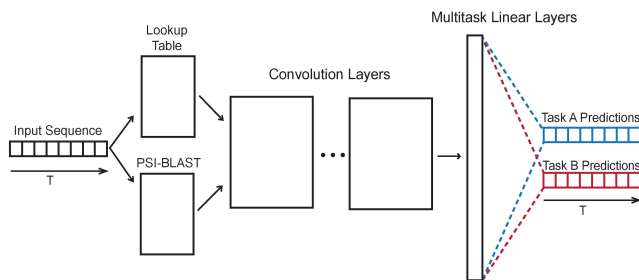


Figure 2: An overview of the deep architecture of our model. Our model accepts an input protein sequence of length  $T$ , which is fed through the network to generate per-position predictions of length  $T$  for several tasks.

This novel multilayer shift-and-stitch technique makes it feasible to train a CNN end-to-end and generate dense per-position protein property prediction. This technique allows us to use convolution and maxpooling layers to label sequences of arbitrary length.

MUST can also be extended to train sequences in minibatches if needed, though the operations will be slightly more complicated. However, we found minibatches not useful, because each amino acid is a training example, and each sequence already contains many amino acids. Additionally, sequences are generally of different lengths, which make implementation of minibatches harder.

### End-to-end Architecture

In this section we describe the end-to-end model structure of the MUST-CNN and how we are able to train it to make fully dense per-position predictions.

The input into the network is a one-hot encoding of an amino acid base pair sequence and the PSI-BLAST position specific scoring matrix (PSSM), which is described in more detail in section Experiments subsection Feature. Dropout is applied to the amino acid input and then fed through a Lookup Table, similar to (Collobert et al. 2011), to construct an embedding representation for each amino acid. Then, the features from the amino acid embeddings are joined directly with the PSSM matrices along the feature dimension and fed into the deep convolutional network.

To apply the shift-and-stitch technique, we shift the amino acid sequences according to the amount of pooling in each layer. Then, we pass every shift through each layer as described above, and stitch the results together after all convolutional layers. This creates a deep embedding for every amino acid in our sequence. Most previous methods use windowing to label the center amino acid. In our model, we can run the whole sequence through the model instead of each window at a time. This allows us to take advantage of the speed of convolution operations and use much larger models.

We use a multitask construction similar to (Qi et al. 2012), where we pass the deep embedding from the convolution layers into several linear fully connected layers which classify the protein sequence into each separate task. This as-

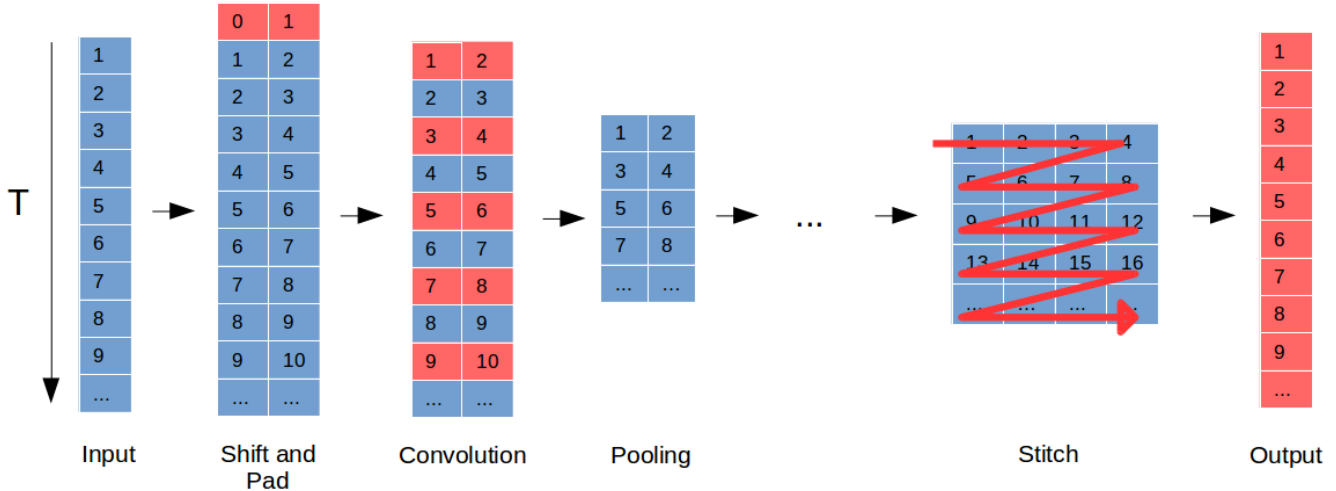


Figure 3: Shift-and-stitch allows us to tag every element of an input even though maxpooling downsamples inputs. By zero padding each sequence correctly, we can join them along the batch dimension and process different shifts at the same time. This technique generalizes arbitrarily to any number of layers, and we can stitch together the result by rearranging and reshaping the resultant tensor, making computation very efficient.

sumes a linear relationship between the deep embedding of a protein chain and the properties predicted. In order for us to classify the outputs of the network for task  $\tau \in \mathcal{T}$ , into class  $c \in C_\tau$  for sequence  $s \in S$ , we apply the softmax operator on the outputs  $f_{t,\tau,c,s}$  of the subclassifiers for task  $\tau$  at position  $t = 1, \dots, T$ . Given the parameters of the network  $\theta$ , this gives us a conditional probability of class  $c$ :

$$p_\tau(c \in C_\tau | f_{t,\tau,s}, \theta) = \frac{e^{f_{t,\tau,c,s}}}{\sum_{c \in C_\tau} e^{f_{t,\tau,c,s}}}$$

The parameters of the network are trained end-to-end by minimizing the negative log-likelihood function over the training set, summing over all tasks and all elements in the sequence:

$$L(\theta) = - \sum_{s \in S} \sum_{\tau \in \mathcal{T}} \sum_{t=1}^T \ln p_\tau(c_{correct} | f_{t,\tau,s}, \theta)$$

where  $c_{correct}$  is the correct label of the amino acid. The minimization of the loss function is obtained via the stochastic gradient descent (SGD) algorithm with momentum, where we update the parameters after every sequence. After the initial multitask model is trained, we take the top layers and each task-specific subclassifier and *fine-tune* the models by initializing their weights at the weights learned by the multitask model and training only on each specific task with  $\frac{1}{10}$  of the original learning rate. Regularization is achieved via dropout (Srivastava et al. 2014).

All models are implemented using the Torch7 framework (Collobert, Kavukcuoglu, and Farabet 2011).

## Connecting to Previous Studies

MUST-CNN is closely related to three previous models: OverFeat (Sermanet et al. 2013), Generative Stochastic networks (GSNs) (Zhou and Troyanskaya 2014), and Conditional Neural Fields (CNFs) (Wang et al. 2011).

CNFs are equivalent to a Conditional Random Field (CRF) with a convolutional feature extractor. As far as we know, the authors implement a windowed version using MLP networks. Their model, although able to consider the entire sequence due to the use of a CRF, is unable to build deeper representations of models. Our model uses multiple convolutional layers and multitasking to classify each amino acid into one of a few classes across multiple tasks. Our models are much deeper, and hence can learn more efficient representations for complex dependencies.

The GSN is similar to a Restricted Boltzmann Machine with interconnections between the hidden states. Training requires a settling algorithm similar to finding the stationary distribution of a Markov chain. Although this technique allows for a model that considers the entire protein sequence, it is less well understood. Convolution layers have the advantage of being used more often in industry (See Related Works), and being well understood. Additionally, a fully feedforward model is almost certainly faster than a model that requires a distribution to converge, though (Zhou and Troyanskaya 2014) did not state training or testing time in their paper.

OverFeat is the most closely related, though it works on images instead of sequence based classification. The pipeline of OverFeat takes in images and classifies them densely to detect objects at every patch. Then the bounding boxes for the objects are combined into a single bounding box, which is used to localize the object. MUST-CNN

is a one dimensional classification algorithm, which takes in the protein sequence surrounding an amino acid and returns a dense property prediction of each amino acid. However, since object localization does not need to be done on every bounding box, OverFeat only uses shift-and-stitch on the last layer for a small resolution improvement. We do fully end-to-end shift-and-stitch, which is difficult on the image domain due to the quadratic increase in calculation time.

## Experiments

### Feature

The features that we use are (1) individual amino acids and (2) PSI-BLAST information (Altschul et al. 1997) of a protein sequence. Each amino acid  $a \in A$ , where  $A$  is the dictionary of amino acids, is coded as a one-hot vector in  $\mathbb{R}^{|A|}$ . That is, the encoding  $x$  of the  $i$ -th amino acid has  $x_i = 1$  and  $x_{j \neq i} = 0$ . PSI-BLAST generates a PSSM of size  $T \times 20$  for a  $T$  lengthed sequence, where a higher score represents a higher likelihood of the  $i^{\text{th}}$  amino acid replacing the current one in other species. Generally, two amino acids that are interchangeable in the PSSM indicates that they are also interchangeable in the protein without significantly modifying the functionality of the protein. The PSI-BLAST profiles were generated in the same way as the original authors in each of the datasets (Qi et al. 2012; Zhou and Troyanskaya 2014).

### Data

We used two large protein property datasets in our experiments. The train, validation and test splits are given in Table 1. The two datasets we use are as follows:

*4prot* Derived from (Qi et al. 2012), we use a train-validation-test split where the model is trained on the training set, selected via validation set results, and best results reported by testing on the test set.

*CullPDB* Derived from (Zhou and Troyanskaya 2014), we choose to use the CullPDB dataset where sequences with  $> 25\%$  identity with the CB513 dataset was removed. The train and validation sets are derived from CullPDB while the test set is CB513 in order to compare results with (Kaae Sønderby and Winther 2014; Wang et al. 2011).

### Tasks

Both datasets were formatted to have the same multitask representation. These are the four classification tasks we tested our system on:

*dssp* The 8 class secondary structure prediction task from the dssp database (Touw et al. 2015). The class labels are H = alpha helix, B = residue in isolated beta bridge, E = extended strand, G = 3-helix, I = 5-helix, T = hydrogen bonded turn, S = bend, L = loop.

*ssp* A collapsed version of the 8 class prediction task, since many protein secondary structure prediction algorithms use a 3 class approach instead of the 8-class approach given in dssp.  $\{H, G\} \rightarrow H = \text{Helix}$ ,  $\{B, E\} \rightarrow B = \text{Beta sheet}$ , and  $\{I, S, T, L\} \rightarrow C = \text{Coil}$

Datasets	Number of	train	validation	test
4prot	Protein chains	7076	2359	2359
	Amino Acids	1500k	509k	506k
CullPDB & CB513	Protein chains	4427	1107	513
	Amino Acids	949k	235k	85k

Table 1: Size of datasets. We do a 60-20-20 split between training, test, and validation datasets on 4prot, but a 80-20-0 split on CullPDB, since we are testing on CB513.

	MUST-CNN small	MUST-CNN large
Convolution Layers	3	3
Hidden units	189	1024
Convolution Size	9	5
Maxpooling Size	2	2
Input Dropout	.35	.1
Dropout	0	{.5, .3}
Nonlinearity	ReLU	ReLU
Learning Rate	0.0148	0.01
Momentum	0.9	0.9

Table 2: Model parameters for all models. The parameters on the small model were discovered via Bayesian Optimization, while the parameters on the large model were discovered using grid search assisted manual tuning. The dropout on the large network was 0.5 on CullPDB, but 0.3 on 4prot, adjusted based on the difference between training and validation error. All models were trained for 50 iterations.

*sar* Relative solvent accessibility. Given the most solvent accessible amino acid in the protein has  $x \text{ \AA}$  of accessible surface area, we label other amino acids as solvent accessible if they have greater than  $0.15x \text{ \AA}$  of accessible surface area.

*sa* Absolute solvent accessibility. Defined as the amino acid having more than  $0.15 \text{ \AA}$  of accessible surface area.

### Training

*Model Selection (Small model)* We use Bayesian Optimization (Snoek, Larochelle, and Adams 2012) to find the optimal model. This is done using the Spearmint package (Snoek 2015). We ran Bayesian Optimization for one week to find the optimal parameters for the small model.

*Model Selection (Large model)* The large model was found using a combination of grid search and manual tuning. The specific architectures we found is detailed in Table 2. Bayesian Optimization could not be used because large models were too slow to train.

After training of the joint model, we also fine-tuned the model by considering each individual task and kickstarting the training from the models learned in the joint model. That is, we started training a model whose parameters were the same as the multitask model, but the loss function only included one specific task. The loss function for task  $\tau$ , sequence  $s$  indexed from  $t = 1, \dots, T$  is then

$$L(\theta) = - \sum_{s \in S} \sum_{t=1}^T \ln p_{\tau}(c_{correct} | f_{t,\tau,s}, \theta)$$

This result is labeled as fine-tune in tables 2 and 3 We use the validation set during the finetuning to find the best dropout value, but then we include the validation set in the retraining set. Dropout generally ensures that early stopping is not needed, so including the validation set should improve the accuracy of our model. We fine-tune at a learning rate of  $\frac{1}{10}$  of the joint model learning rate.

*Time* Training of the small model takes 6 hours, while training of the large model takes one day. Since testing the fine-tuned models involve passing the data through four separate models, while testing the multitask model involves doing all at the same time, it takes longer to test on the fine-tuned model. Nevertheless, we were able to handle testing speeds of over a million amino acids in under 2 seconds.

*Hardware* In order to speed up computation, we utilize the parallel architecture of the GPU, which is especially useful for convolutional models which do many parallel computations. All training and testing uses a Tesla C2050 GPU unit.

## Results

During model selection, we discovered that our model is very robust to model parameters. Most combinations of parameter tweaks inside the optimal learning rate give a less than 1% improvement in average accuracy. By using max-pooling with shift-and-stitch in our model our average accuracy improved by almost 0.5% with barely any computational slowdown.

Our results on the 4prot dataset are detailed in Table 3. The small model we found via Bayesian Optimization has approximately as many parameters as previous state-of-the-art models, but we see that it outperformed the network created by (Qi et al. 2012) on all tasks. Fine-tuning on individual models is necessary for good performance. This implies that it may perhaps be easier to build an MLP subclassifier for each task, instead of assuming linearity. Training jointly on the large model already beats (Qi et al. 2012), but fine-tuning increases the accuracy dramatically. Additionally, the testing time is reported in milliseconds per million amino acids. We see that the small models can test fairly quickly, while the fine-tuned large models have a  $2.5\times$  slowdown. We are the first to report precise training and testing times for a model on protein property prediction.

A detailed listing of precision-recall scores for 4prot is given in Table 4. We see the expected pattern of lower frequencies having a lower F1 score, since unbalanced datasets are harder to classify. Precision is very stable, while recall dramatically lowers according to the frequency of labels. This suggests that our model picked up on several key properties of labels with few training examples, but missed many. More training data is one way to solve this issue.

Our results on the CullPDB dataset and comparisons with existing state-of-the-art algorithms is detailed in ta-

Task	Qi et al.	Conv small	fine-tuned	Conv large	fine-tuned
dssp (8)	68.2	67.0	70.6	69.5	<b>76.7</b>
ssp (3)	81.7	80.6	84.0	82.5	<b>89.6</b>
sar (2)	81.1	79.0	81.2	80.2	<b>84.9</b>
saa (2)	82.6	80.9	82.9	82.0	<b>86.1</b>
Test time	596k*	379	587	553	1597

Table 3:  $Q_c$  accuracy on different architectures of model on 4prot dataset. The number in parenthesis behind the task determines  $c$ , the number of classes in each task. Testing time is given for all tasks simultaneously in *milliseconds per million amino acids*. (\*) Test time was not detailed in referenced paper, so their algorithm was implemented and tested on the CPU.

Per-task Label	Recall	Precision	F1	Frequency
<b>dssp</b>				
H	.967	.878	.920	.328
E	.924	.821	.869	.206
L	.748	.645	.693	.211
T	.564	.623	.592	.113
S	.254	.621	.360	.095
G	.363	.655	.467	.035
B	.049	.797	.093	.012
I	0	0	0	.0002
<b>ssp</b>				
C	.875	.881	.878	.418
H	.936	.919	.928	.364
E	.868	.884	.876	.218
<b>sar</b>				
Inaccessible	.874	.838	.856	.512
Accessible	.823	.861	.842	.488
<b>saa</b>				
Inaccessible	.901	.888	.894	.650
Accessible	.789	.810	.799	.350

Table 4: Recall, precision, and F1 scores for 4prot dataset. Class I for dssp does not occur often enough for our model to learn labelings.

ble 5. We do 1% better than the previous published best, despite using a dramatically simpler algorithm. Testing on the CB513 dataset allows a direct comparison to how previous methods perform. We do not achieve a dramatically higher accuracy rate as we do on 4prot. We suspect that filtering non-homologous protein sequences decreases possible accuracy, since we are essentially demanding a margin of difference between the data distributions for the training and testing samples. It may not be possible to predict protein properties accurately using a statistical method if non-homologous protein sequences were filtered from the training set.

## Discussion

We have described a multilayer shift-and-stitch convolutional architecture for sequence prediction. We use ideas from the image classification domain to train a deep convolutional network on per-position sequence labeling. We are the first to use multilayer shift-and-stitch on protein se-

Model	$Q_8$
CNF (Wang et al. 2011)	.649
GSN (Zhou and Troyanskaya 2014)	.664
LSTM (Kaae Sønderby and Winther 2014)	.674
MUST-CNN (Ours)	<b>.684</b>

Table 5:  $Q_8$  accuracy training on the CullPDB dataset and testing on CB513. Testing takes around the same time as for the 4prot dataset. We use the same architecture as MUST-CNN large, detailed in table 2.

quences to generate per-position results. Shift-and-stitch is a trick to quickly compute convolutional network scores on every single window of a sequence at the same time, but the fixed window sizes of the convolutional network still remains. Surprisingly, we achieve better results than whole sequence-based approaches like the GSN, LSTM, and CNF models used in previous papers (see Table 5). We believe this is because the speed of our model allows us to train models with far higher capacity. We show that the architecturally simpler MUST-CNN does as well or better than more complex approaches.

In our experiments, the same network works very well on two different large datasets of protein property prediction, in which we only changed the amount of dropout regularization. This suggests that our model is very robust and can produce good results without much manual tuning once we find a good starting set of hyperparameters. More generally, our technique should work on arbitrary per-position sequence tagging tasks, such as part of speech tagging and semantic role labeling.

Additionally, our model can make predictions for a million amino acids in under 2 seconds. Although the main speed bottleneck of protein property prediction is obtaining the PSI-BLAST features, the speed of our model can be useful on other sequence prediction tasks where feature extraction is not the bottleneck.

Future work can incorporate techniques such as the fully convolutional network (Long, Shelhamer, and Darrell 2014) to further speed up and reduce the parameter set of the model. Another direction is to continue along the lines of LSTMs and GSNs and try to better model the long range interactions of the protein sequences.

## References

Altschul, S. F.; Madden, T. L.; Schäffer, A. A.; Zhang, J.; Zhang, Z.; Miller, W.; and Lipman, D. J. 1997. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Research* 25(17):3389–3402.

Bahdanau, D.; Cho, K.; and Bengio, Y. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.

Collobert, R., and Weston, J. 2008. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, 160–167. ACM.

Collobert, R.; Weston, J.; Bottou, L.; Karlen, M.; Kavukcuoglu, K.; and Kuksa, P. 2011. Natural language processing (almost) from scratch. *The Journal of Machine Learning Research* 12:2493–2537.

Collobert, R.; Kavukcuoglu, K.; and Farabet, C. 2011. Torch7: A matlab-like environment for machine learning. In *BigLearn, NIPS Workshop*.

Cuff, J. A., and Barton, G. J. 2000. Application of multiple sequence alignment profiles to improve protein secondary structure prediction. *Proteins: Structure, Function, and Bioinformatics* 40(3):502–511.

Drozdetskiy, A.; Cole, C.; Procter, J.; and Barton, G. J. 2015. JPred4: a protein secondary structure prediction server. *Nucleic Acids Research* gkv332.

Glorot, X.; Bordes, A.; and Bengio, Y. 2011. Deep sparse rectifier neural networks. In *International Conference on Artificial Intelligence and Statistics*, 315–323.

He, K.; Zhang, X.; Ren, S.; and Sun, J. 2015. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *arXiv preprint arXiv:1502.01852*.

Jones, D. T. 1999. Protein secondary structure prediction based on position-specific scoring matrices. *Journal of Molecular Biology* 292(2):195–202.

Kaae Sønderby, S., and Winther, O. 2014. Protein Secondary Structure Prediction with Long Short Term Memory Networks. *ArXiv e-prints*.

Kim, Y. 2014. Convolutional Neural Networks for Sentence Classification. *arXiv:1408.5882 [cs]*. arXiv: 1408.5882.

Lecun, Y.; Bottou, L.; Bengio, Y.; and Haffner, P. 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE* 86(11):2278–2324.

Long, J.; Shelhamer, E.; and Darrell, T. 2014. Fully convolutional networks for semantic segmentation. *arXiv preprint arXiv:1411.4038*.

Magnan, C. N., and Baldi, P. 2014. SSpro/ACCpro 5: almost perfect prediction of protein secondary structure and relative solvent accessibility using profiles, machine learning and structural similarity. *Bioinformatics (Oxford, England)* 30(18):2592–2597.

Pinheiro, P. H. O., and Collobert, R. 2013. Recurrent Convolutional Neural Networks for Scene Parsing. *arXiv:1306.2795 [cs]*. arXiv: 1306.2795.

Qi, Y.; Oja, M.; Weston, J.; and Noble, W. S. 2012. A unified multitask architecture for predicting local protein properties. *PLoS one* 7(3):e32235.

Scherer, D.; Müller, A.; and Behnke, S. 2010. Evaluation of pooling operations in convolutional architectures for object recognition. In *Artificial Neural Networks–ICANN 2010*. Springer. 92–101.

Sermanet, P.; Eigen, D.; Zhang, X.; Mathieu, M.; Fergus, R.; and LeCun, Y. 2013. Overfeat: Integrated recognition, localization and detection using convolutional networks. *arXiv preprint arXiv:1312.6229*.

Snoek, J.; Larochelle, H.; and Adams, R. P. 2012. Practical Bayesian optimization of machine learning algorithms. In

*Advances in neural information processing systems*, 2951–2959.

Snoek, J. 2015. Spearmint Bayesian optimization codebase. <https://github.com/HIPS/Spearmint>.

Srivastava, N.; Hinton, G.; Krizhevsky, A.; Sutskever, I.; and Salakhutdinov, R. 2014. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research* 15(1):1929–1958.

Sutskever, I.; Vinyals, O.; and Le, Q. V. 2014. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, 3104–3112.

Szegedy, C.; Liu, W.; Jia, Y.; Sermanet, P.; Reed, S.; Anguelov, D.; Erhan, D.; Vanhoucke, V.; and Rabinovich, A. 2014. Going Deeper with Convolutions. *arXiv:1409.4842 [cs]*. arXiv: 1409.4842.

Touw, W. G.; Baakman, C.; Black, J.; te Beek, T. A. H.; Krieger, E.; Joosten, R. P.; and Vriend, G. 2015. A series of PDB-related databanks for everyday needs. *Nucleic Acids Research* 43(D1):D364–D368.

Wang, Z.; Zhao, F.; Peng, J.; and Xu, J. 2011. Protein 8-class secondary structure prediction using conditional neural fields. *Proteomics* 11(19):3786–3792.

Zaremba, W.; Sutskever, I.; and Vinyals, O. 2014. Recurrent Neural Network Regularization. *arXiv:1409.2329 [cs]*. arXiv: 1409.2329.

Zhou, J., and Troyanskaya, O. G. 2014. Deep supervised and convolutional generative stochastic network for protein secondary structure prediction. *arXiv preprint arXiv:1403.1347*.