

## Ex04: Comparison of different frameworks - ('Implementing a neural part-of-speech tagger')

**Models:** DyNet 2.0.3, PyTorch 0.4.1, Tensorflow 1.9.0 & **Model Design (same for all approaches):** word embeddings initialized using GloVe, one-layer bi-directional LSTM, matrix multiplication to produce scores for tags

Difference	DyNet	PyTorch	TensorFlow
<i>Random Seed Initialization: Guaranteeing reproducible results</i>	<p>random seed needs to be set in <code>dynet_config</code> before the actual import of DyNet</p> <p>+ setting seed ensures reproducible results; random number generator is initialized with the same starting value</p>	<p><code>manual_seed(0)</code> sets the random seed after import of PyTorch</p> <p>+ setting seed ensures reproducible results; random number generator is initialized with the same starting value</p>	<p>no random seed setting done in this project =&gt; but <code>tf.set_random_seed()</code> exists in TensorFlow 1.9</p> <p>- defining a random seed at the start is important to produce reproducible results. - if no random seed is set TensorFlow automatically picks a random seed (doesn't guarantee reproducibility since it changes)</p>
<i>Consistent Size of Input Vectors</i>	<p>no special padding symbol required but kept in code for consistency =&gt; done by DyNet automatically</p> <p>+ practical that vectors are automatically padded</p>	<p>Special PAD symbol introduced</p> <p>- manual padding required not bad but additional effort, user needs to be aware of this</p>	<p>Special PAD symbol introduced</p> <p>- manual padding required not bad but additional effort, user needs to be aware of this</p>
<i>Model Definition</i>	<p>Model parameters (e.g., layers, optimizer) defined first, and the computations are only defined when required (during training i.e. computation graph is built in <code>do_pass</code> function)</p> <p>+ logical definition (1<sup>st</sup> design the model, 2<sup>nd</sup> computations = optimize/train the model) - LSTM initialization parameters are done a lot more explicit than in other models</p>	<p>Model parameters defined as part of constructor of a class, the computations are defined inside functions of the class (e.g., <code>forward()</code> function).</p> <p>+ training = method, classification task = class =&gt; analogy to object-oriented programming + reusing structures - unsure whether the class design has any real benefits</p>	<p>Both the parameters and the computations must be defined at the start. Starts with new computation graph and placeholders for variables fed into graph during computation.</p> <p>+ More definition power (I would assume since every variable is defined manually) - Very verbose: placeholders + all the definitions at the start make it harder to understand at the beginning. More code required to do the same than with other frameworks.</p>
<i>Gradient Clipping</i>	<p>Automatically built-in and activated (disable manually in this implementation)</p> <p>+ ensures that exploding gradients problem is automatically taken care of - automatic activation can also be negative (user needs to be aware of this fact, otherwise he searches for other reasons first)</p>	<p>PyTorch Utilities has two methods: <code>clip_grad_norm</code> and <code>clip_grad_value</code> which would allow to clip the gradients manually.</p> <p>+/- no per default gradient clipping + easier than in TensorFlow - user needs to check gradients if training diverges (optimizer oscillates due to too large gradient steps)</p>	<p>Gradient clipping needs to be performed manually:</p> <ul style="list-style-type: none"> <li>- compute gradients</li> <li>- clip/modify gradients</li> <li>- apply modified gradients</li> </ul> <p>+/- no per default gradient clipping - user needs to check gradients if training diverges (optimizer oscillates due to too large gradient steps)</p>

<b>Difference</b>	<b>DyNet</b>	<b>PyTorch</b>	<b>TensorFlow</b>
<i>LSTM Cell Definition</i>	<p>Requires two LSTM cells to create a bi-directional RNN.</p> <p>+ more definition power (cells could potentially differ) - requires the definition of two separate cells +/- initialization is a bit more effort but potentially more definition freedom (different initializations)</p>	<p>Single LSTM cell is sufficient. Bi-directionality can be specified via an input parameter in the constructor of the LSTM class.</p> <p>+ efficient and convenient for the end-user - less definition power/freedom</p>	<p>To create a bi-directional RNN two separate LSTM cells need to be defined (a forward and a backward LSTM cell)</p> <p>+ more definition power / possibilities for the end-user - for any action the LSTM cells need to be wrapped - more work than PyTorch</p>
<i>Learning Rate Scheduling</i>	<p>Computation manually defined and updated during each epoch of the training</p> <p>+ less boilerplate code than PyTorch, update is defined, directly computed and then reassigned</p>	<p>Scheduler: defined together with the optimizer and the initial learning rate. Defined as a lambda function which receives the current number of epochs at each iteration and then multiplies it by the initial learning rate.</p> <p>+ once defined, method invocation at the start of each epoch - boilerplate code (more lines of code to achieve the same compared to TensorFlow or DyNet)</p>	<p>Computation manually defined and updated during each epoch of the training</p> <p>+ less boilerplate code than PyTorch, update is defined, directly computed and then reassigned</p>
<i>Recurrent Dropout Regularization</i>	<p>Can be set directly on each LSTM (e.g., <code>b_lstm.set_dropouts(0.0)</code> =&gt; disable input and output variational dropout of this cell)</p> <p>+ very simple to define like TensorFlow directly on LSTM cell</p>	<p>PyTorch 0.4 doesn't support recurrent dropout directly =&gt; use toolkit (WeightDrop class)</p> <p>- no support for recurrent dropout, only via separate class</p>	<p>Can be specified via parameter of Dropout wrapper cell for each LSTM cell.</p> <p>+ simple to define directly on LSTM cell like DyNet</p>
<i>Input / Output Dropout</i>	<p>Is applied manually to input (<code>dy.dropout(weights, 1-PROB)</code>) and automatically rescales the outputs</p> <p>+ very practical: only needs to be applied during training since DyNet automatically rescales the outputs by <code>1/PROB</code></p>	<p>Input (Word) and Output (LSTM) dropout are defined in the constructor of the TaggerModel class and applied during training (calls of the forward method)</p> <p>+/- (neutral – not very bad): dropout is manually applied before and after LSTM cell</p>	<p>Basic LSTM cell wrapped in a Dropout Wrapper Cell which applies the dropout to both input and output. Dropout probabilities for input and output can be set individually.</p> <p>- dropout is always applied and cannot be turned off during prediction. <b>Reason:</b> The computation graph cannot be changed. <b>Solution:</b> dropout probability = 1.0 during prediction to keep all values.</p>
<i>Weight Decay</i>	<p>Defined globally in <code>dynet_config</code> (at import time) =&gt; contains parameter <code>weight_decay</code></p> <p>+ allows for weight decay - fixed weight decay (rescales weights with the same factor at each update)</p>	<p>Directly defined inside the optimizer initialization as a constant</p> <p>+ allows for weight decay, defined logically =&gt; in optimizer - fixed weight decay (same decay at each update)</p>	<p>TensorFlow 1.9 doesn't support weight decay yet =&gt; fixed by pull-request in TensorFlow 1.10</p> <p>- no weight decay possible</p>