

Tricks for Training Neural Models

(Some slides by Yoav Goldberg, Graham Neubig)

Optimization Choices

- **Use Gated Recurrent Unit.**
- **Adaptive learning rate.**
 - adaptive optimizers such as Adam ([Kingma14](#)) because they can better handle the complex training dynamics of RNNs
- **Gradient clipping.**
 - Print or plot the gradient norm to see its usual range
 - then scale down gradients that exceeds this range.
 - This prevents spikes in the gradients to mess up the parameters during training.
- **Early Stopping**
- **Stacked recurrent networks.**
 - Recurrent networks need a quadratic number of weights in their layer size.
 - More efficient to stack two or three smaller layers instead of one big one.
 - Sum the outputs of all layers instead of using only last one, similar to a ResNet/ DenseNet.

Normalizations

- **Normalizing the loss.** (To get losses of similar magnitude across datasets)
 - sum the loss terms along the sequence and divide them by the maximum seq length.
 - This makes it easier to reuse hyper parameters between experiments.
 - The loss should be averaged across the batch.
- **Batch Normalization.** fixes the means and variances of each layer's inputs (Lofe & Szegedy 15)
- **Layer Normalization.** Adding layer normalization (Ba et al 16) to all linear mappings of the recurrent network speeds up learning

Batch Normalization

For a layer of the network with d -dimensional input, $\mathbf{x} = (x^{(1)}, \dots, x^{(d)})$, each dimension of its input is then normalized (i.e. re-centered and re-scaled) separately,

$$\hat{x}_i^{(k)} = \frac{x_i^{(k)} - \mu_B^{(k)}}{\sqrt{\sigma_B^{(k)^2 + \epsilon}}, \text{ where } k \in [1, d] \text{ and } i \in [1, m]; \mu_B^{(k)} \text{ and } \sigma_B^{(k)^2} \text{ are the per-}$$

dimension mean and variance, respectively.

$$y_i^{(k)} = \gamma^{(k)} \hat{x}_i^{(k)} + \beta^{(k)}, \quad \text{adaptive bias and gain}$$

where the parameters $\gamma^{(k)}$ and $\beta^{(k)}$ are subsequently learned in the optimization process.

$$BN_{\gamma^{(k)}, \beta^{(k)}} : \mathbf{x}_{1\dots m}^{(k)} \rightarrow \mathbf{y}_{1\dots m}^{(k)} \text{ called the Batch Normalizing transform.}$$

BN for hidden-hidden transitions is not beneficial; BN for input-hidden transitions can speed up learning!

Layer Normalization

Given inputs x over a minibatch of size m , $B = \{x_1, x_2, \dots, x_m\}$, each sample x_i contains K elements, i.e. the length of flatten x_i is K , by applying transformation of your inputs using some learned parameters γ and β , the outputs could be expressed as $B' = \{y_1, y_2, \dots, y_m\}$, where $y_i = \text{LN}_{\gamma, \beta}(x_i)$.

$$\hat{x}_{i,k} = \frac{x_{i,k} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

$$\mu_i = \frac{1}{K} \sum_{k=1}^K x_{i,k}$$
$$\sigma_i^2 = \frac{1}{K} \sum_{k=1}^K (x_{i,k} - \mu_i)^2$$

$$y_i = \gamma \hat{x}_i + \beta \equiv \text{LN}_{\gamma, \beta}(x_i)$$

LN can be done for each instance separately – hence it does same computation at training/test time. can work with B=1

Model Parameters (RNN)

- **Learned initial state.**

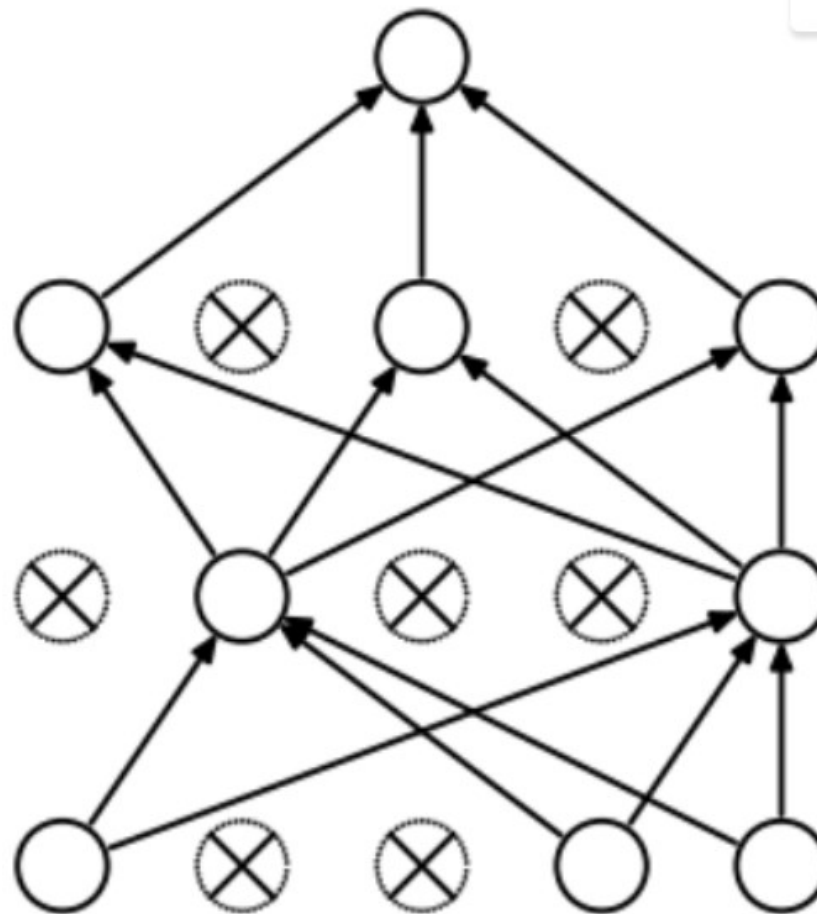
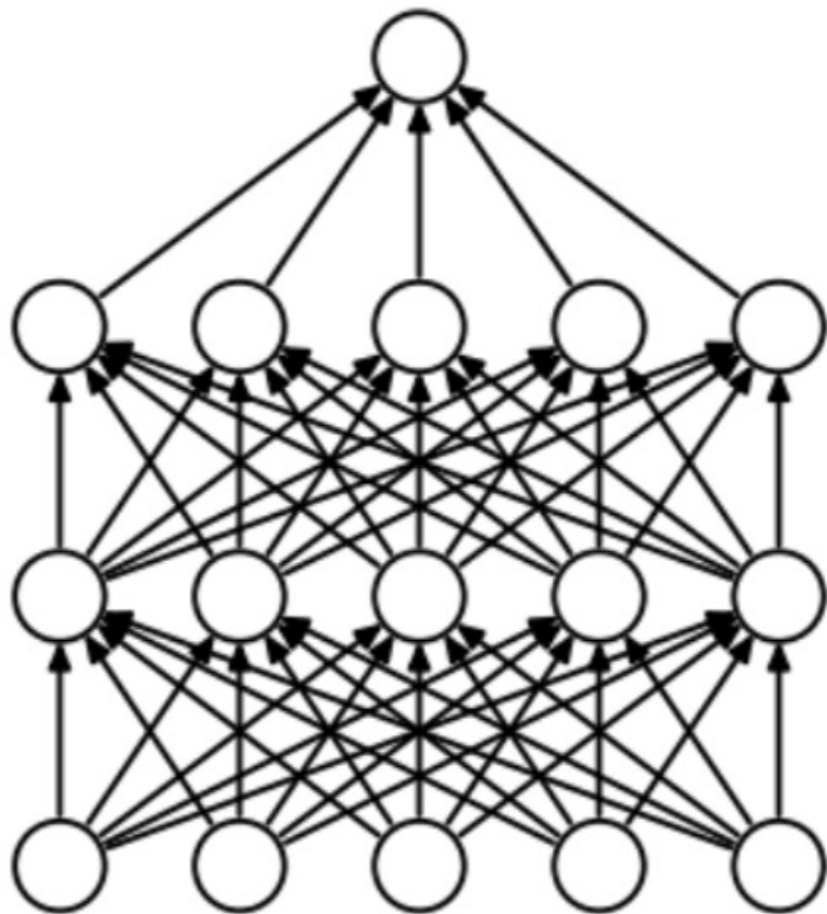
- Initializing the hidden state as zeros → large loss initially
- Training the initial state as a variable can improve performance as described in <https://r2rt.com/non-zero-initial-states-for-recurrent-neural-networks.html>

- **Forget gate bias.**

- It can take a while for a RNN to learn to remember information
- Initialize biases for LSTM's forget gate to 1 to remember more by default.
- Similarly, initialize biases for GRU's reset gate to -1.

- **Regularization.** If your model is overfitting, use dropout

Dropout (Srivastava et al 2014)



Observations on Dropout

Dropout forces a neural network to learn more robust features that are useful in conjunction with many different random subsets of the other neurons.

Dropout roughly doubles the number of iterations required to converge. However, training time for each epoch is less.

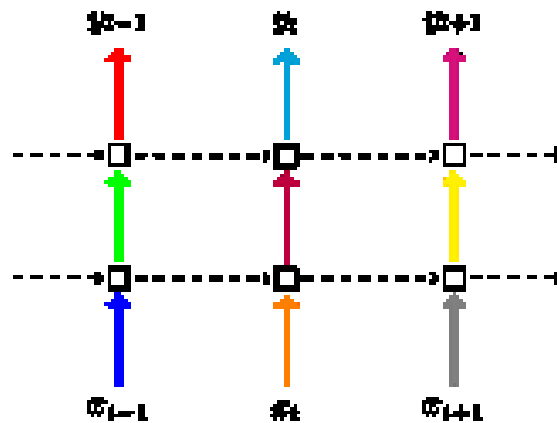
With H hidden units, each of which can be dropped, we have 2^H possible models. In testing phase, the entire network is considered and each activation is reduced by a factor p .

DropConnect

- Wan et al. (2013)
 - Instead of dropping nodes, drop edges (weights)
- Generalization of dropout

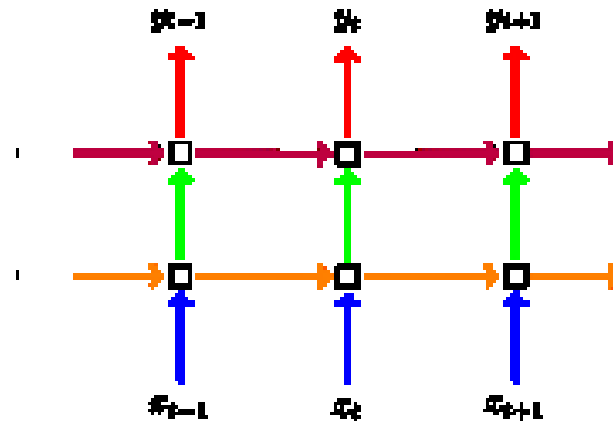
Dropout in RNNs

- Still an open question how to perform well.
- One suggestion: apply only to feedforward part of RNN (Zaremba et al 14)



Dropout in RNNs

- Still an open question how to perform well.
- Yarin Gal's Variational Dropout (Gal & Ghahramani 2015):



uses the same dropout mask at each time step, including the recurrent layers (colours representing dropout masks, solid lines representing dropout, dashed lines representing standard connections with no dropout).

<https://medium.com/@bingobee01/a-review-of-dropout-as-applied-to-rnns-72e79ecd5b7b>

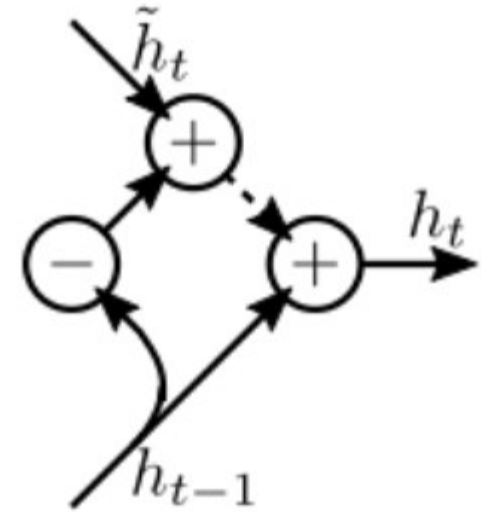
Recurrent Dropout

- (Semeniuta et al 2016) $\mathbf{c}_t = \mathbf{f}_t * \mathbf{c}_{t-1} + \mathbf{i}_t * d(\mathbf{g}_t)$

“We demonstrate that recurrent dropout is most effective when applied to hidden state update vectors in LSTMs rather than to hidden states; (ii) we observe an improvement in the network’s performance when our recurrent dropout is coupled with the standard forward dropout, though the extent of this improvement depends on the values of dropout rates”

ZoneOut

- (Krueger et al 2017)



In a variation on the dropout philosophy, Krueger et al. (2017) proposed Zoneout where *“instead of setting some units’ activations to 0 as in dropout, zoneout randomly replaces some units’ activations with their activations from the previous timestep.”* this *“makes it easier for the network to preserve information from previous timesteps going forward, and facilitates, rather than hinders, the flow of gradient information going backward”*

Ensembles

- **Same model, different initialization.**
 - Use cross-validation to determine the best hyperparameters,
 - then train multiple models with the best set of hyperparameters but with different random initialization.
 - Suffers from limited variety
- **Top models discovered during cross-validation.**
 - Use cross-validation to determine the best hyperparameters
 - then pick the top few (e.g., 10) models to form the ensemble.
 - Improves the variety of ensemble but has the danger of including suboptimal models
- **Different checkpoints of a single model.**
 - If training is very expensive
 - limited success in taking different checkpoints of a single network over time (for example after every epoch) and using those to form an ensemble.
 - Clearly, this suffers from some lack of variety, but can still work reasonably well in practice.
 - The advantage of this approach is that is very cheap.

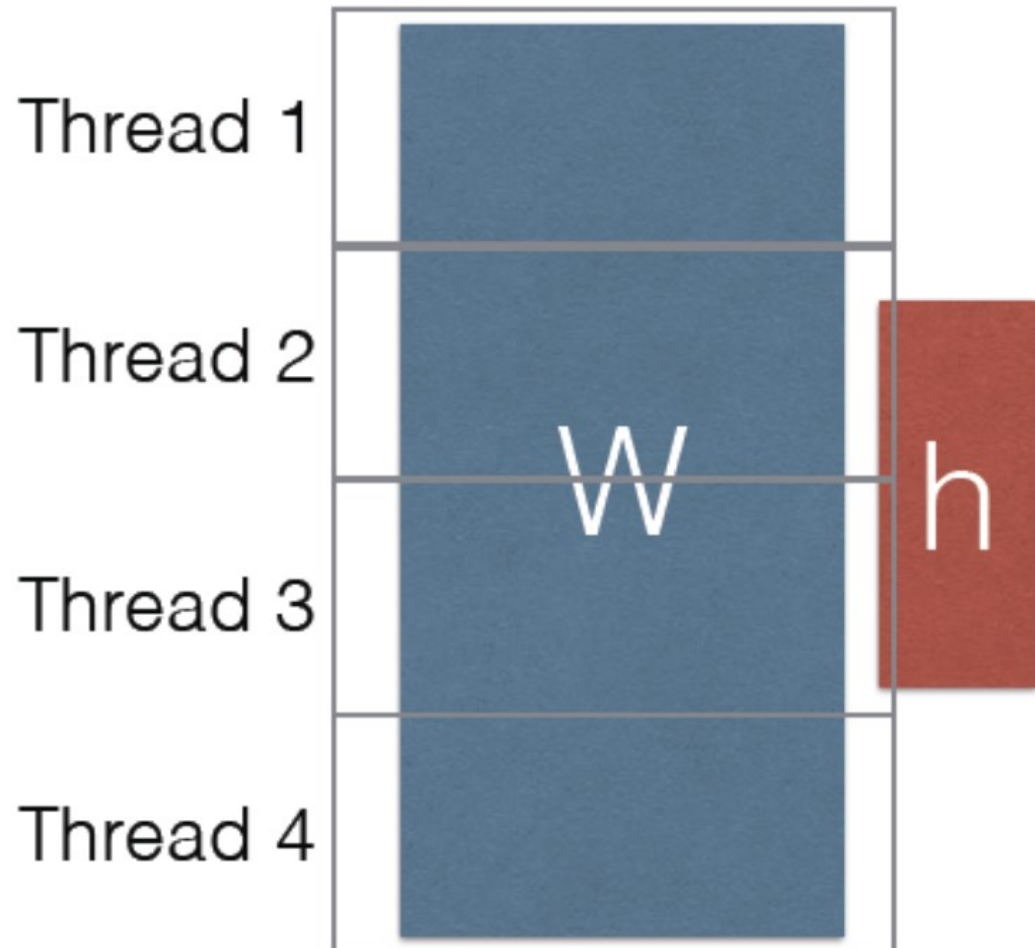
Why are Neural Networks Slow and What Can we Do?

- Big operations, especially for softmaxes over large vocabularies
 - → **Approximate operations or use GPUs**
- GPUs love big operations, but hate doing lots of them
 - → **Reduce the number of operations** through optimized implementations or batching
- Our networks are big, our data sets are big
 - → **Use parallelism** to process many data at once

Parallelism in Computation Graphs

- Three types of parallelism
 - Within-operation parallelism
 - Operation-wise parallelism
 - Example-wise parallelism
- } Model parallelism
- } Data parallelism

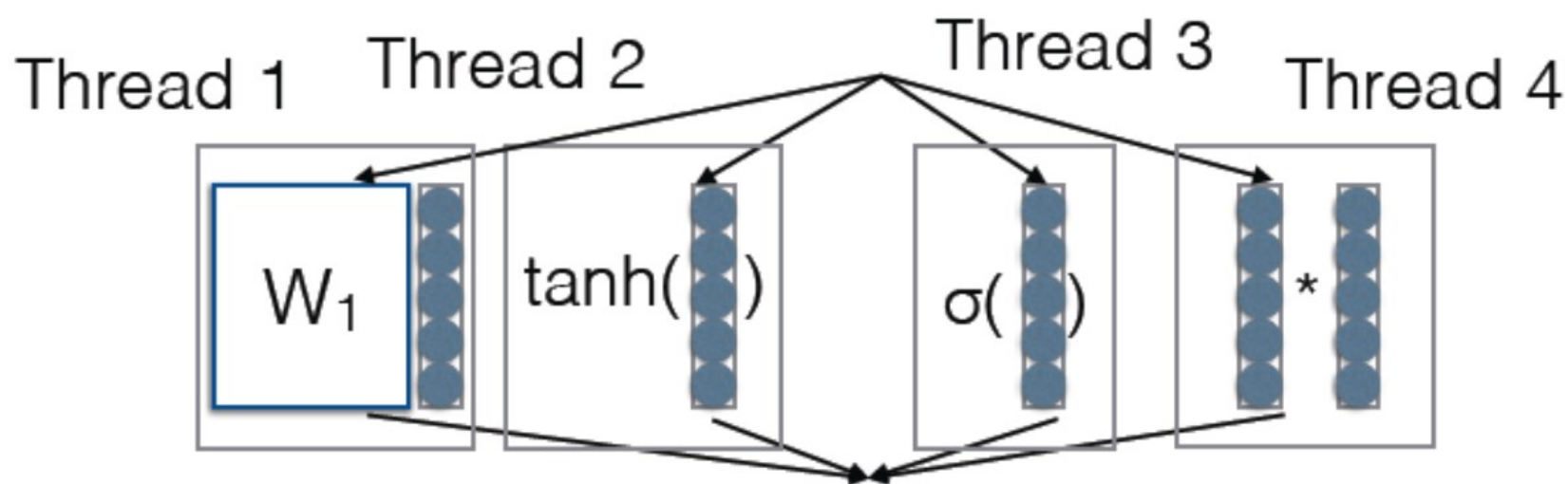
Within-operation Parallelism



- GPUs excel at this!
- Libraries like MKL implement this on CPU, but gains less striking.
- Thread management overhead is counter-productive when operations small.

Operation-wise Parallelism

- Split each operation into a different thread, or different GPU device



- Difficulty:** How do we minimize dependencies and memory movement?

Example-wise Parallelism

- Process each training example in a different thread or machine

this is an example	Thread 1
this is another example	Thread 2
this is the best example	Thread 3
no, i'm the best example	Thread 4

- **Difficulty:** How do we accumulate gradients and keep parameters fresh across machines?

GPUs vs. CPUs

CPU, like a motorcycle



Quick to start, top speed
not shabby

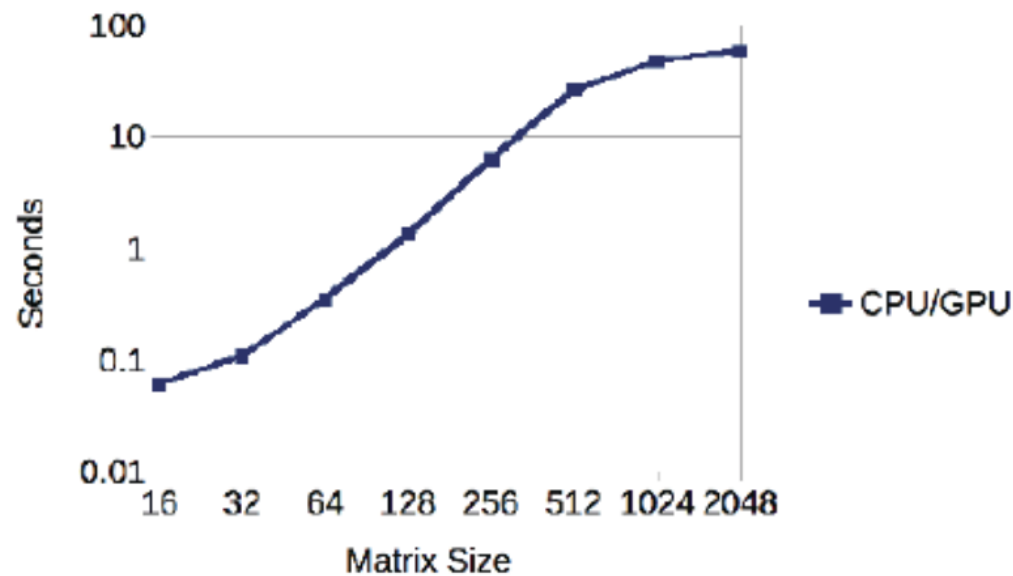
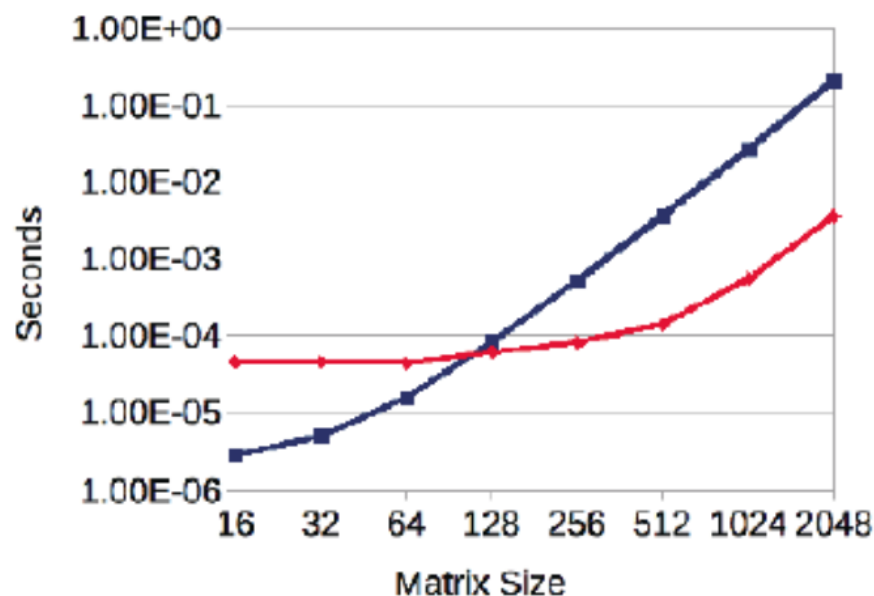
GPU, like an airplane



Takes forever to get off the
ground, but super-fast
once flying

A Simple Example

- How long does a matrix-matrix multiply take?



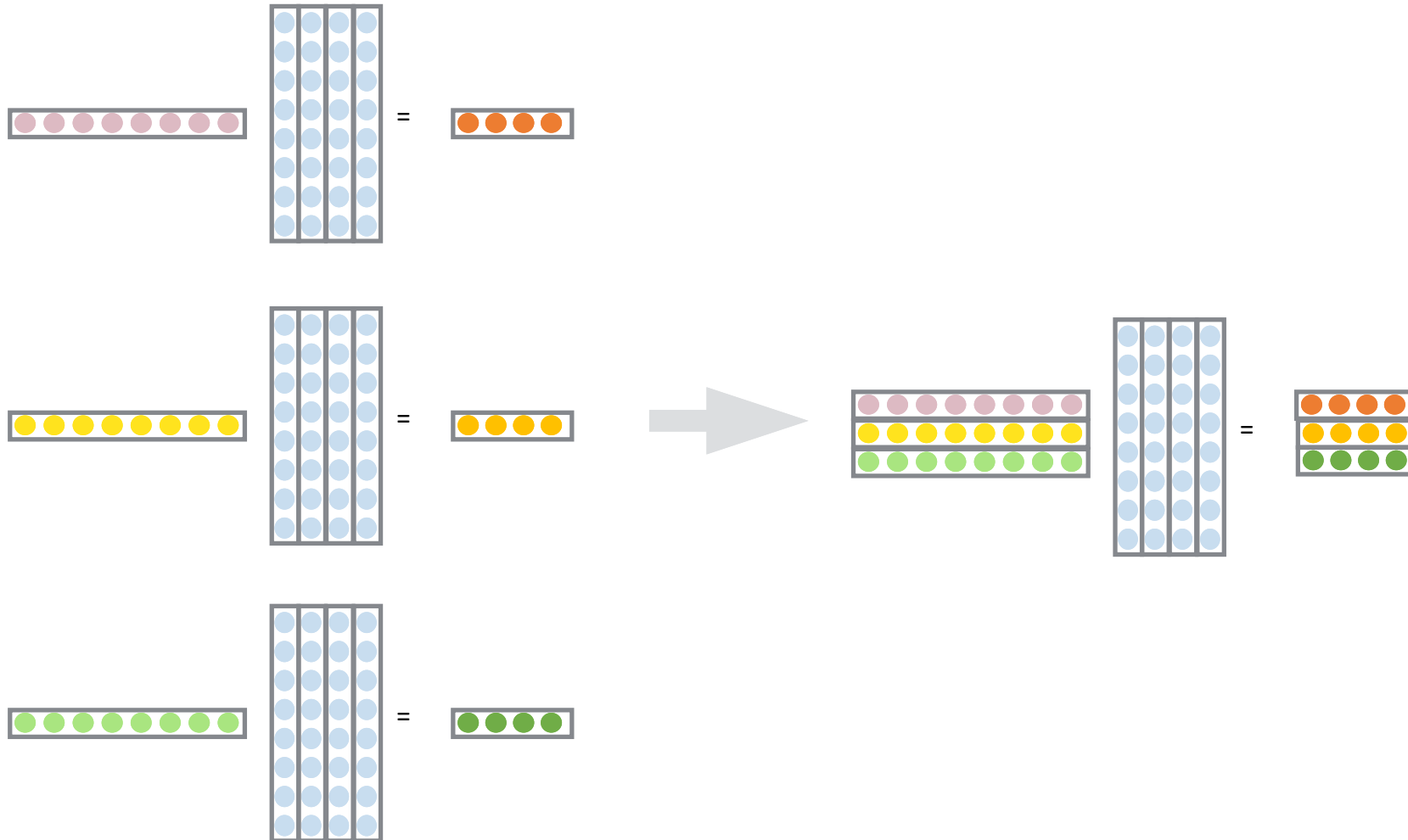
Practically

- Use **CPU for profiling**, it's plenty fast (esp. DyNet) and you can run many more experiments
- For **many applications, CPU is just as fast** or faster than GPU: NLP analysis tasks with small or complicated data/networks
- You see **big gains on GPU when** you have:
 - Very big networks (or softmaxes with no approximation)
 - Do mini-batching
 - Optimize things properly

Batching (in RNNs)

- Most toolkits require a **fixed** computation graph for all examples.
- But RNNs have **different** input lengths. What do we do?
 - Option 1:
Use a tool that does not pose this limitation.
 - Option 2:
Enforce max length + padding for shorter sequences.

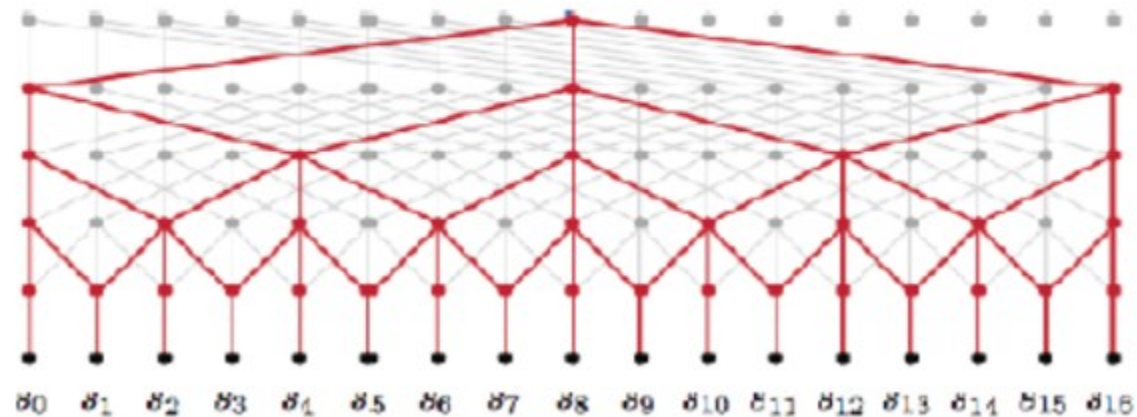
Batching Reminder

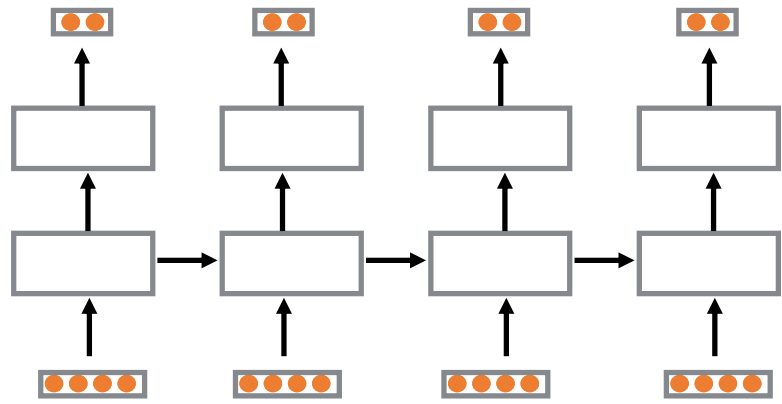
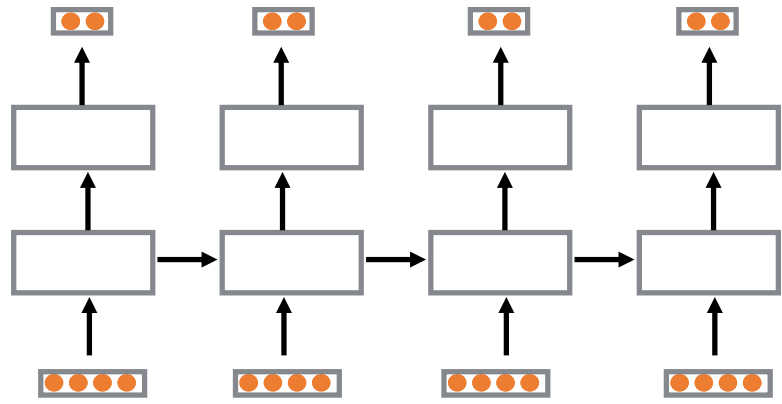
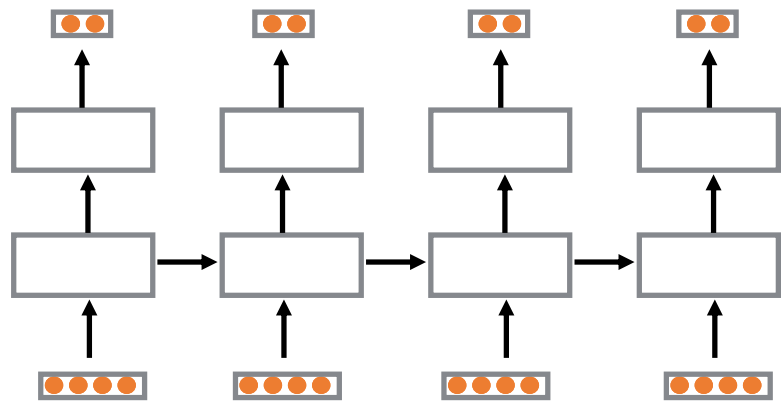


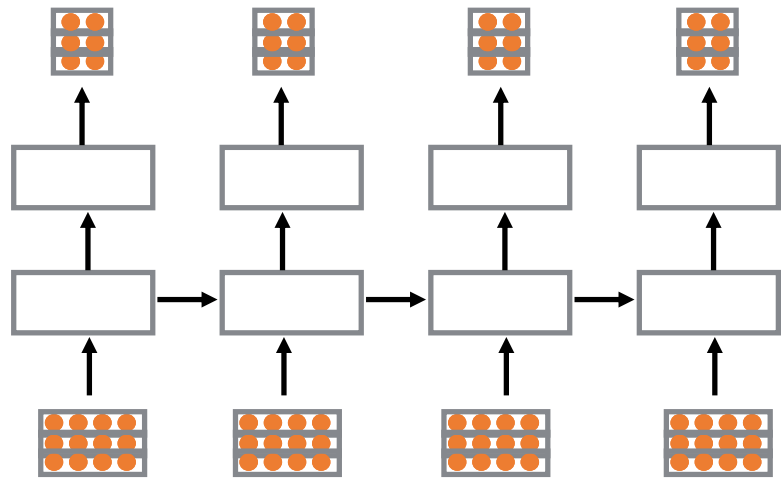
Batching in RNNs

- Sequential in nature, very little parallelism.
- (Compare, e.g., to a Convolutional Network)

Dilated convolutions for capturing context
(Kalchbrenner et al. 2016): single GPU call for entire
sentence!

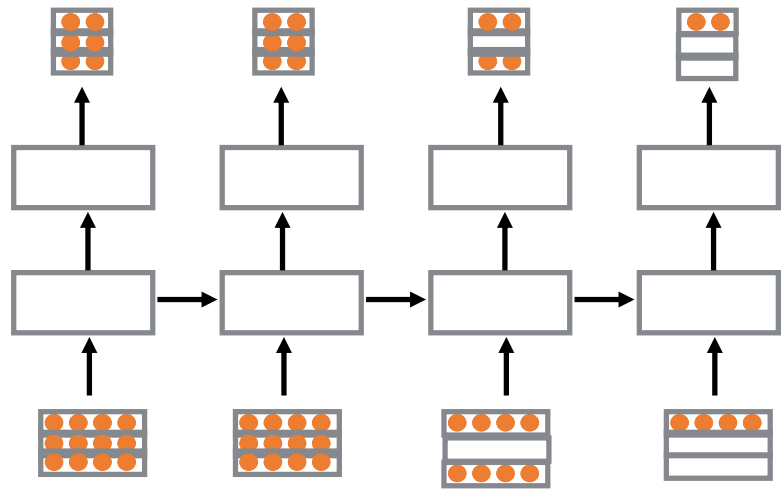


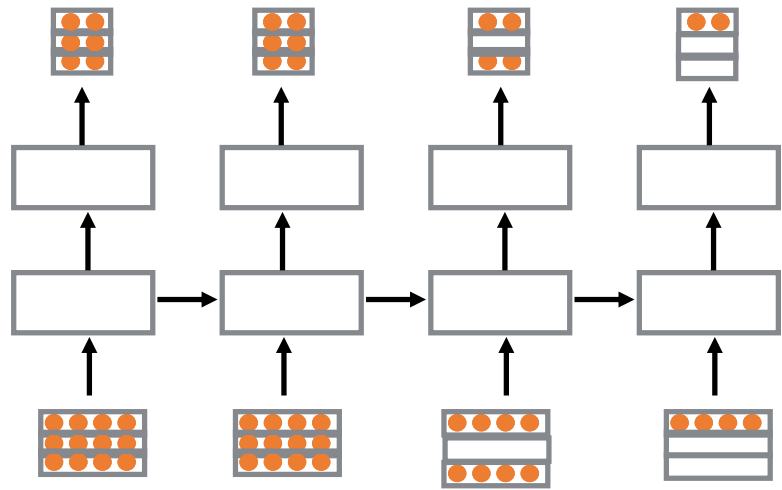




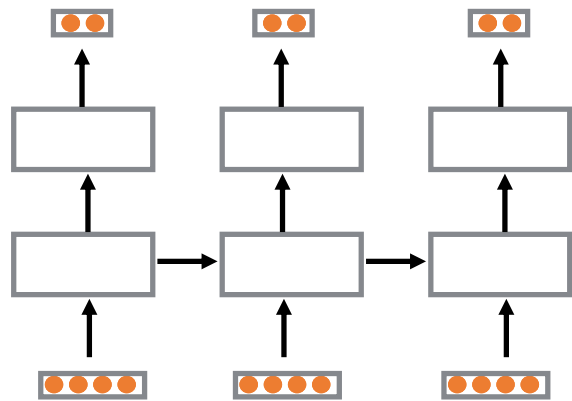
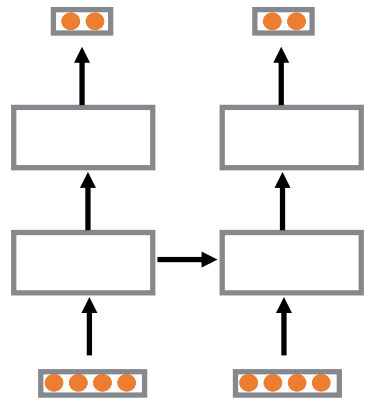
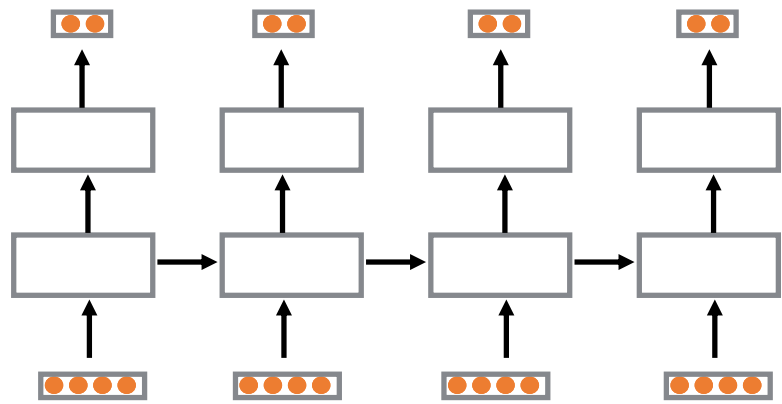
what if the sequences are different lengths?



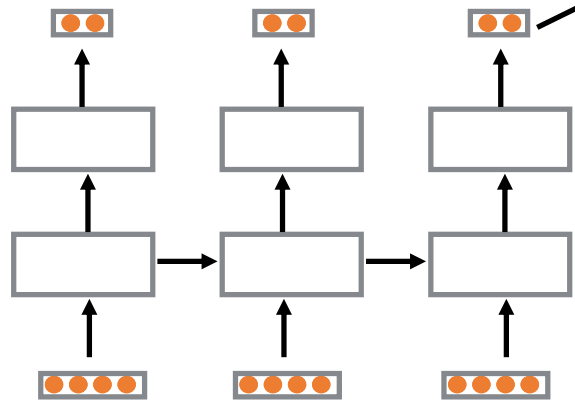
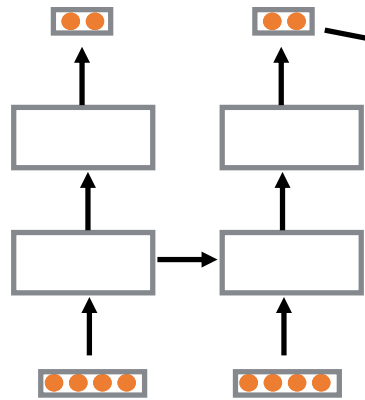
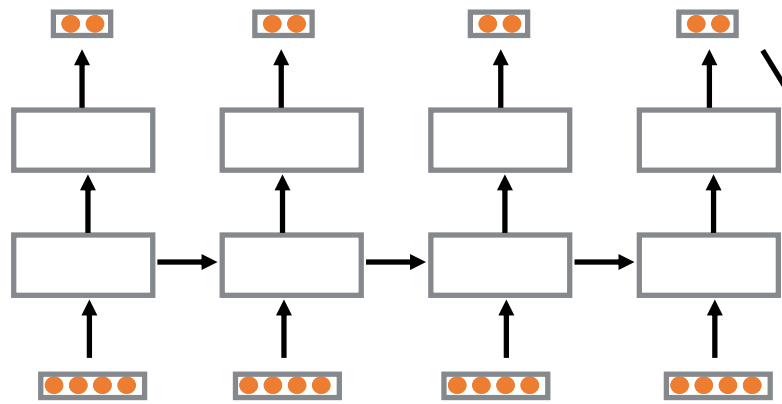




padding

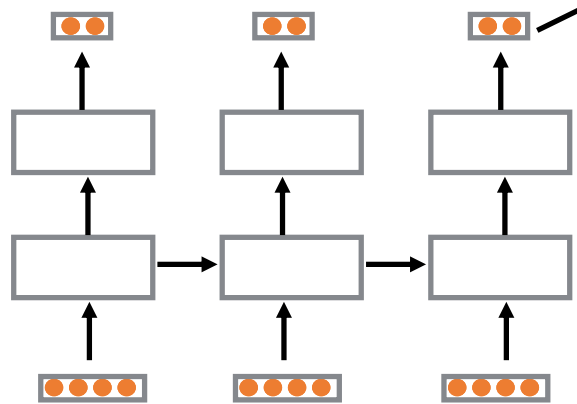
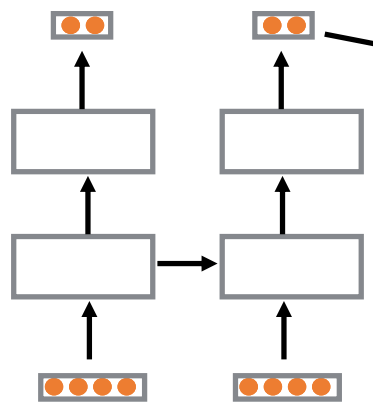
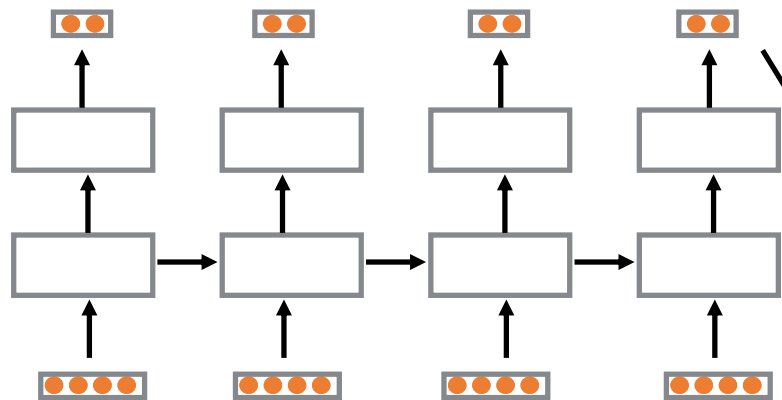


create a separate
network for each
(easy)



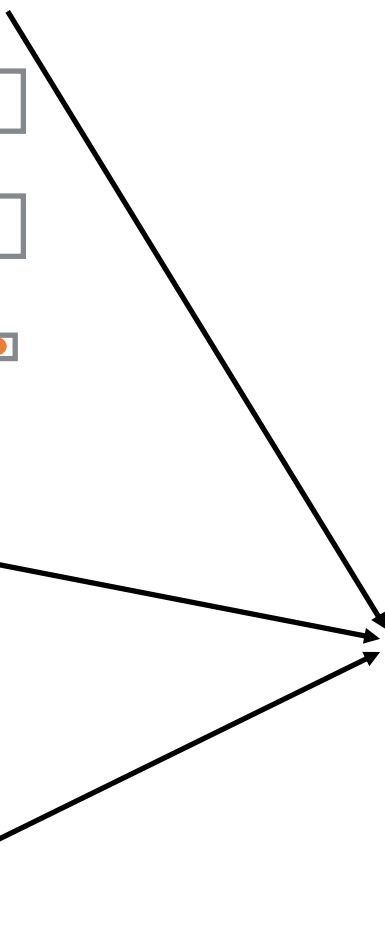
treat them
as a single
graph

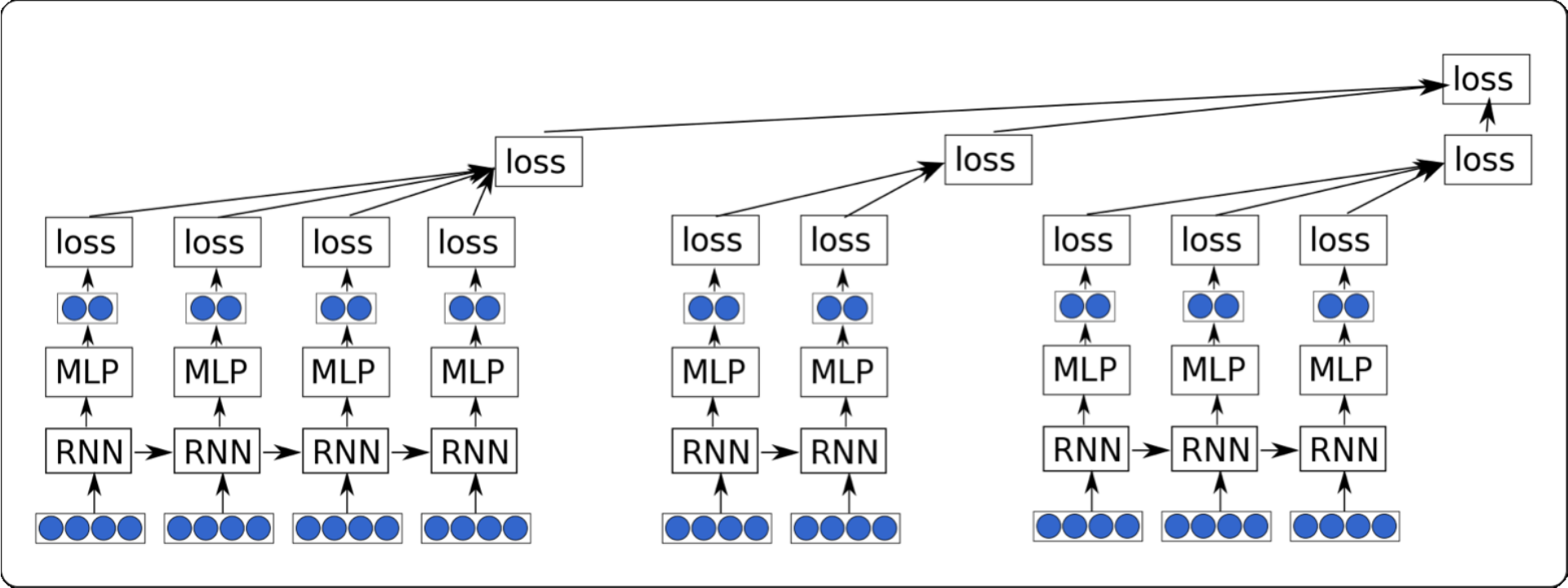




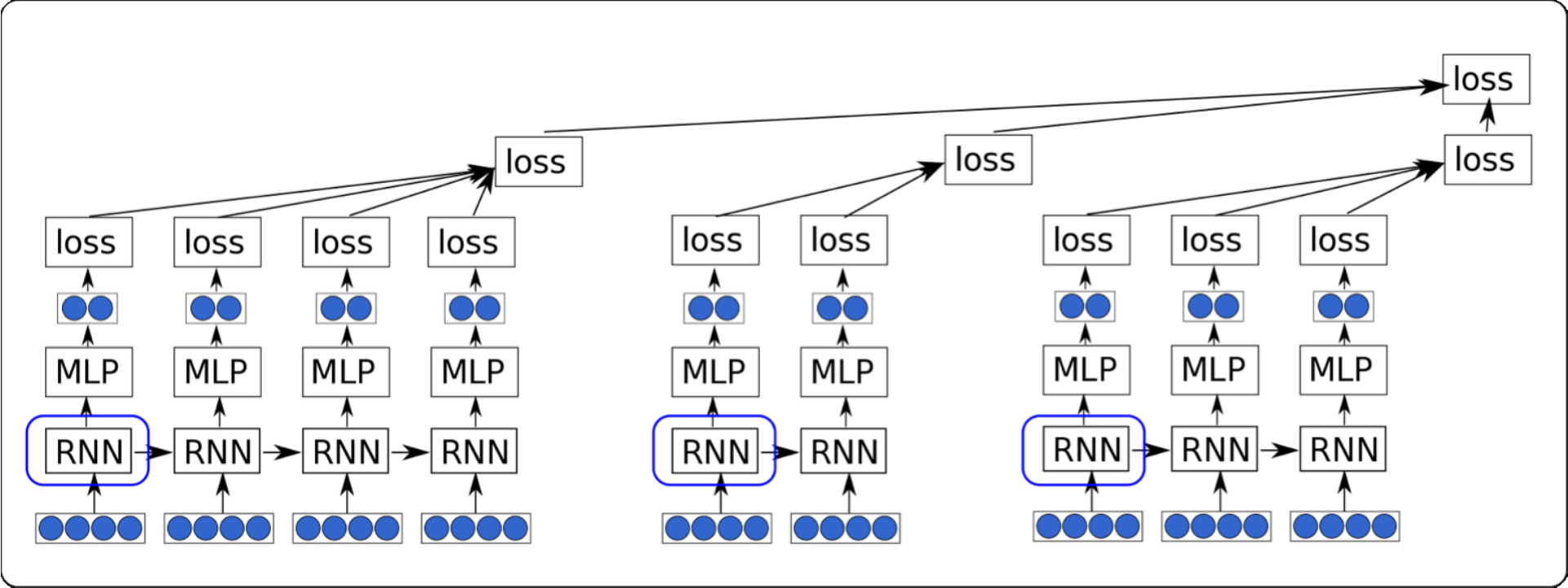
treat them
as a single
graph

Dynet (PyTorch(?)) will identify batching
opportunities for you.

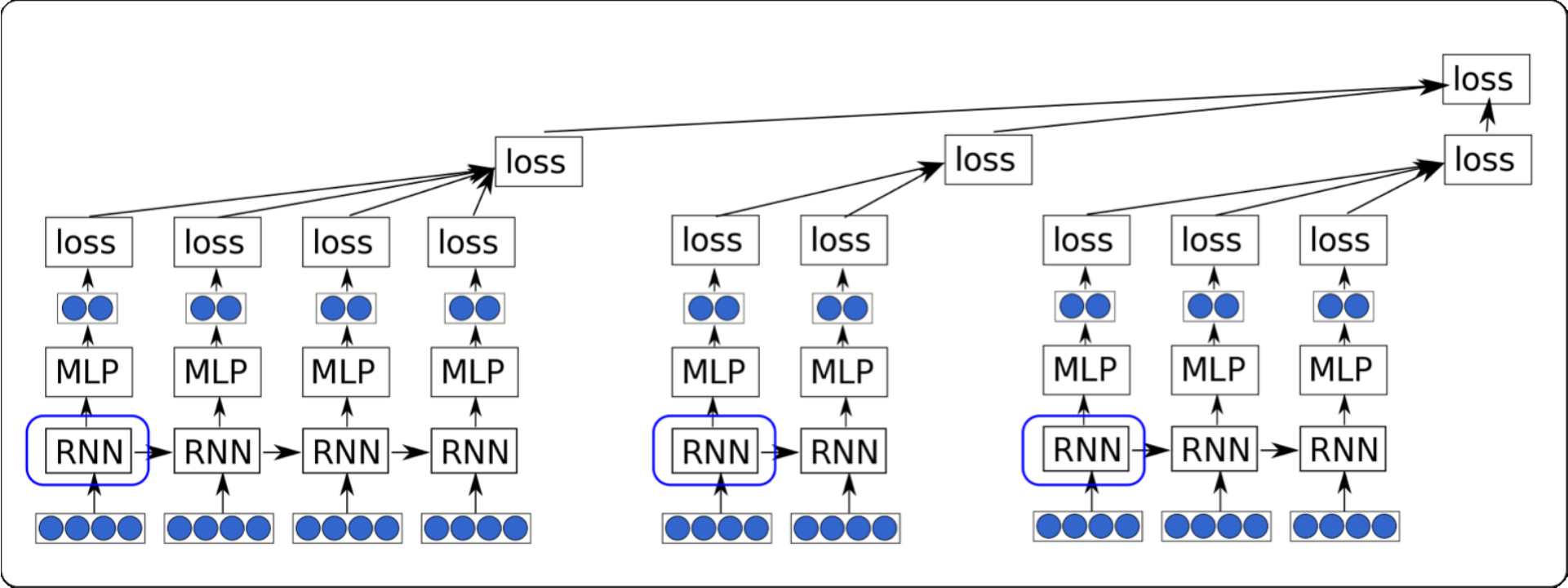


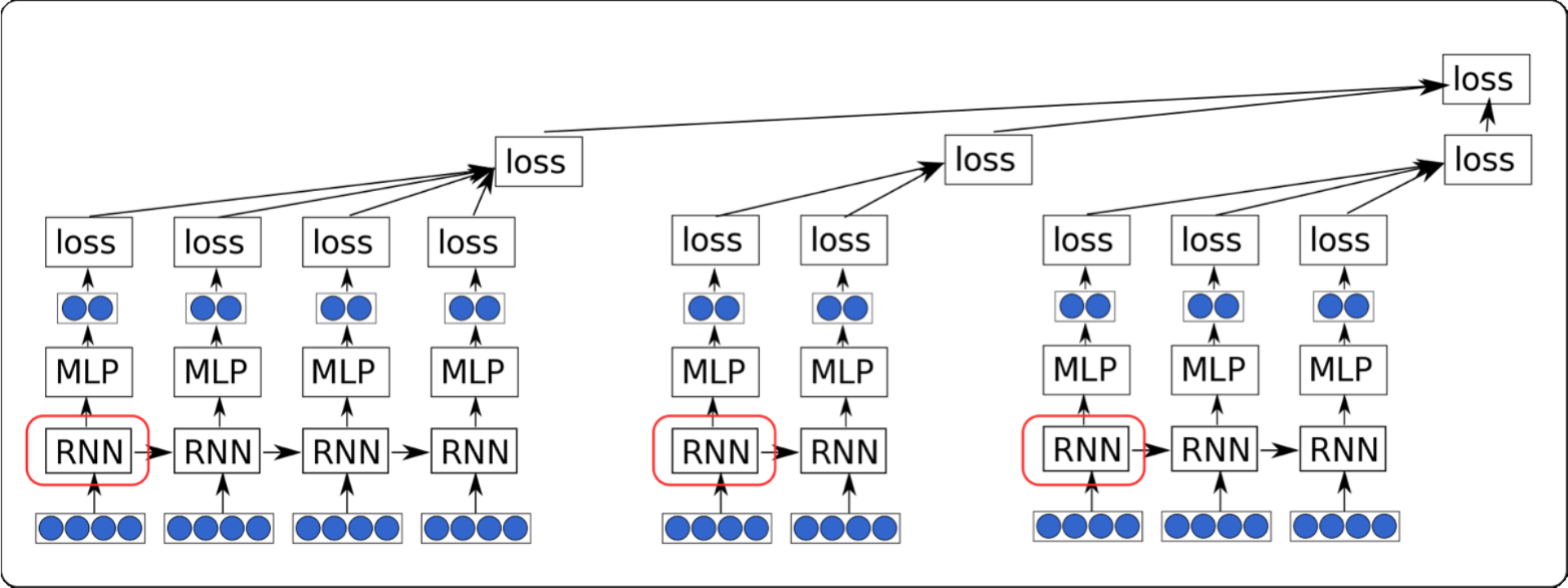


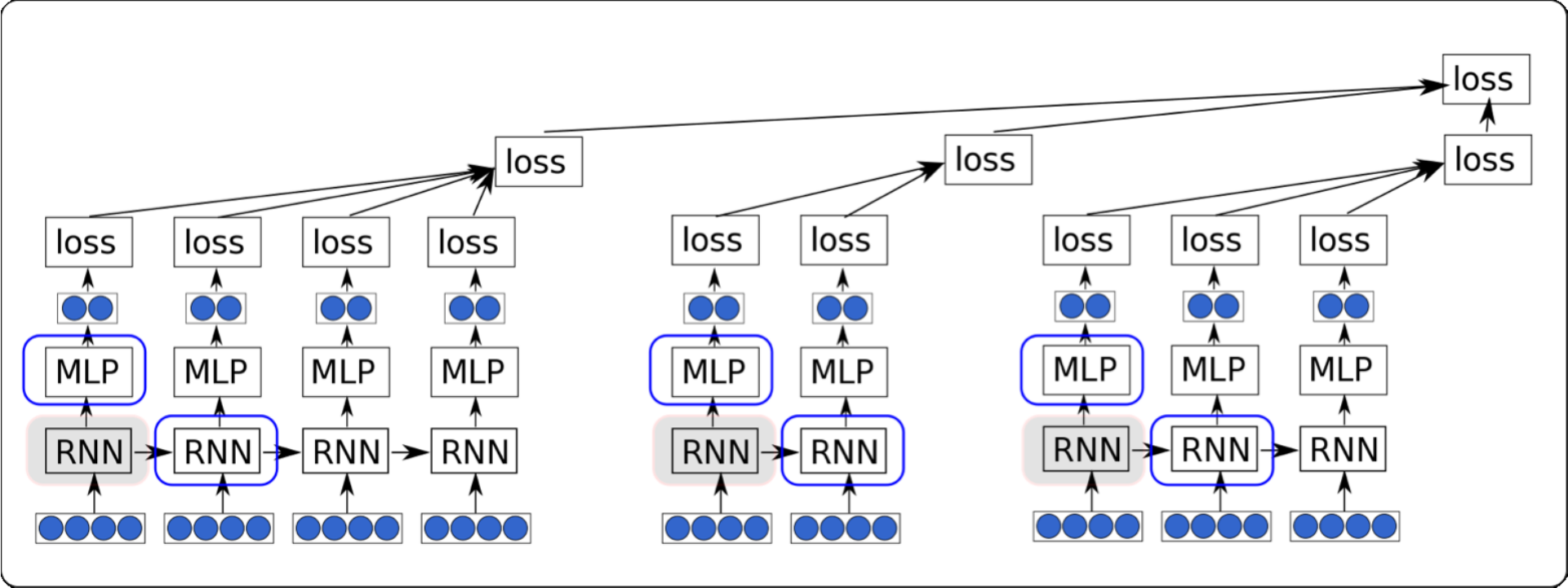
nodes in **blue** are ready
to be executed

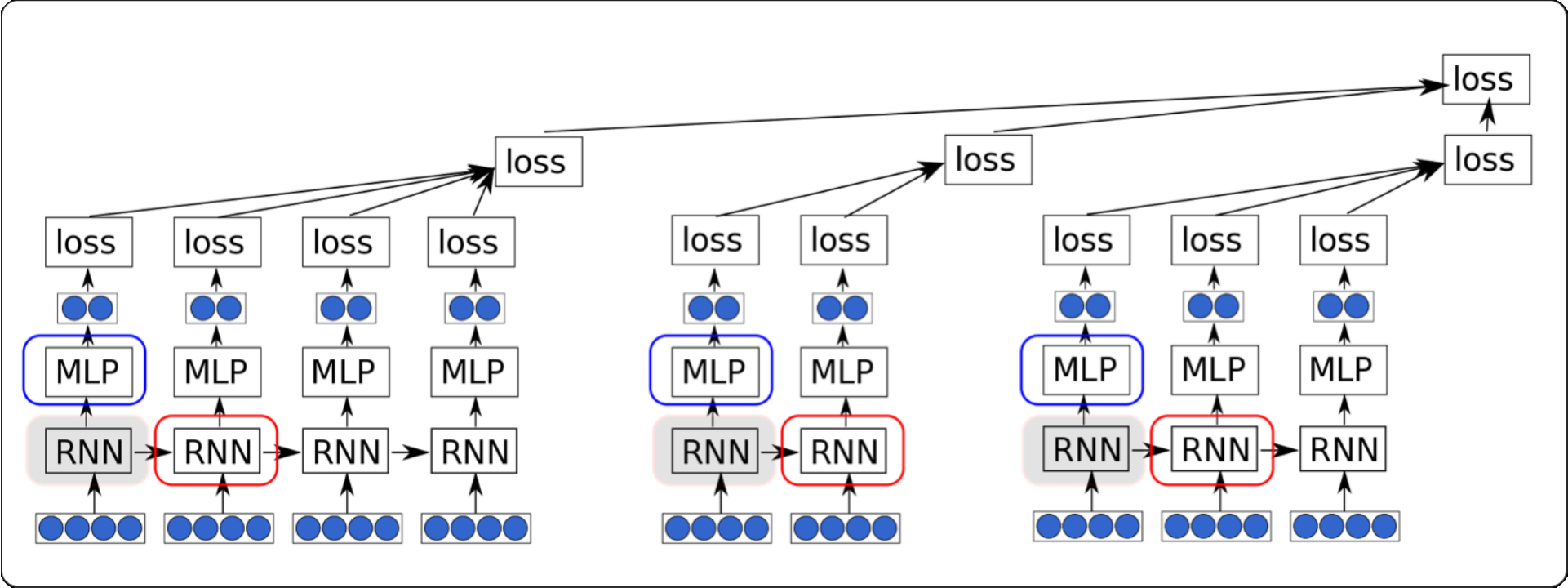


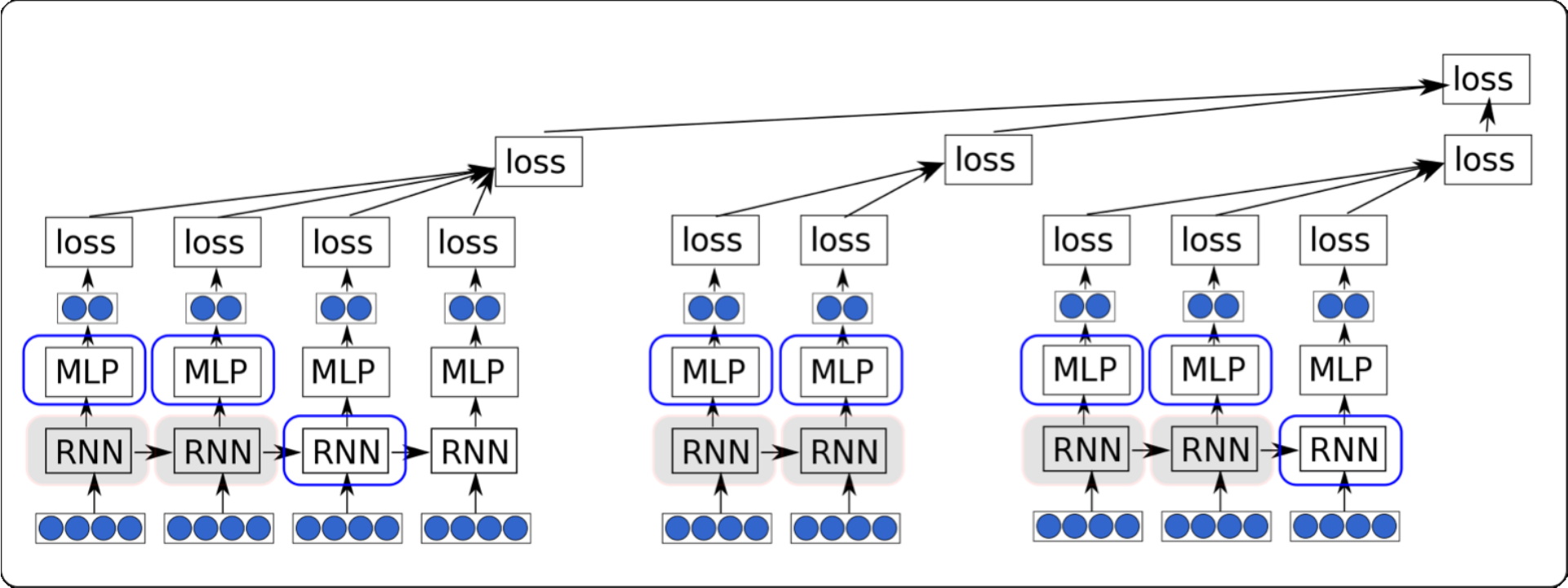
nodes in **red** will be executed
using batch operations

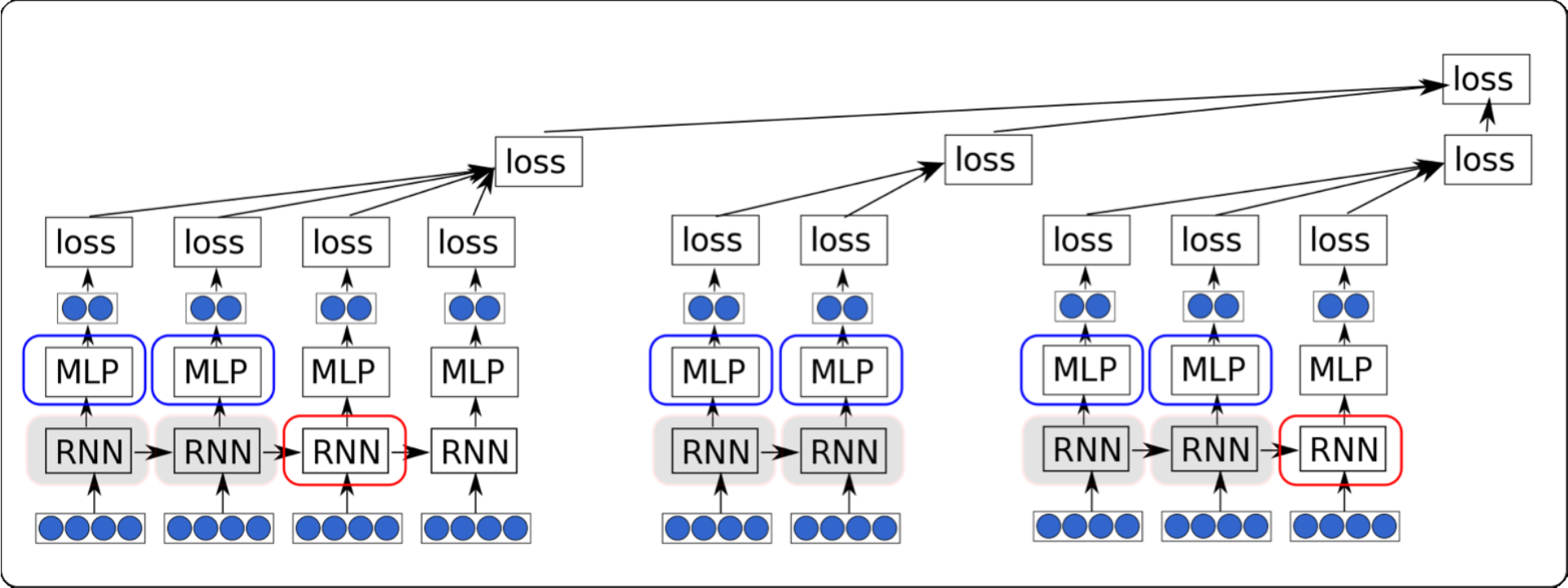


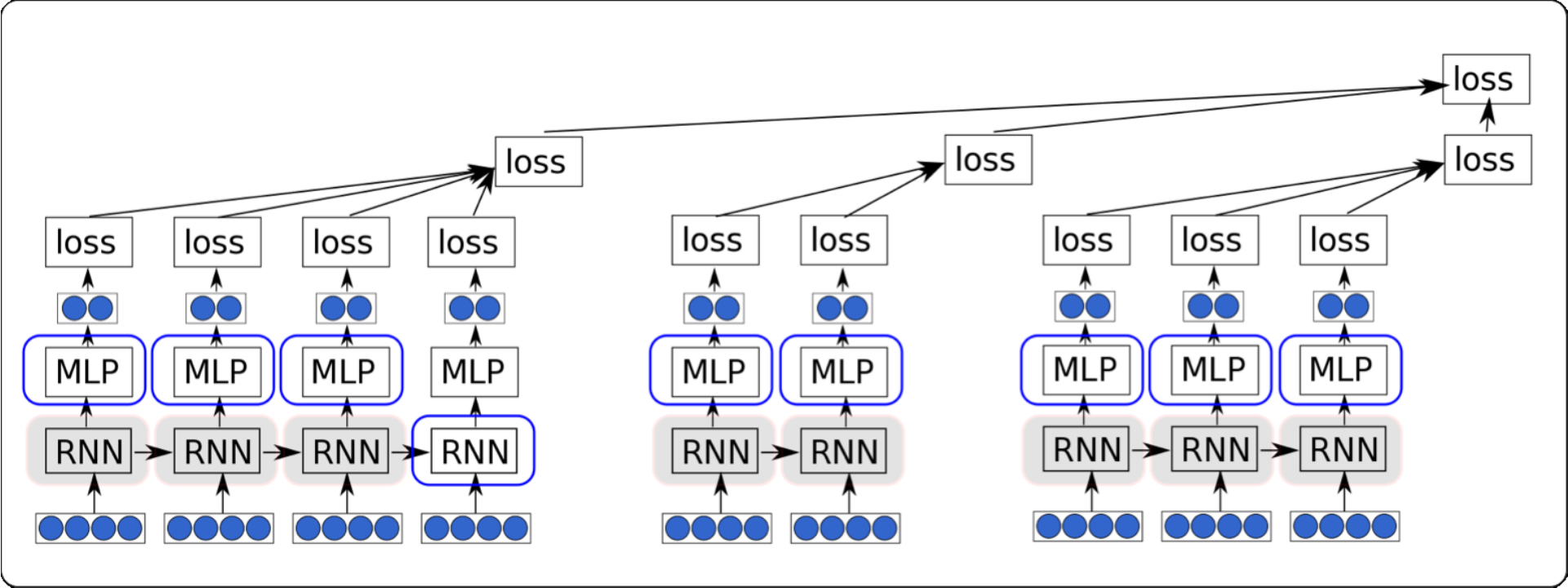


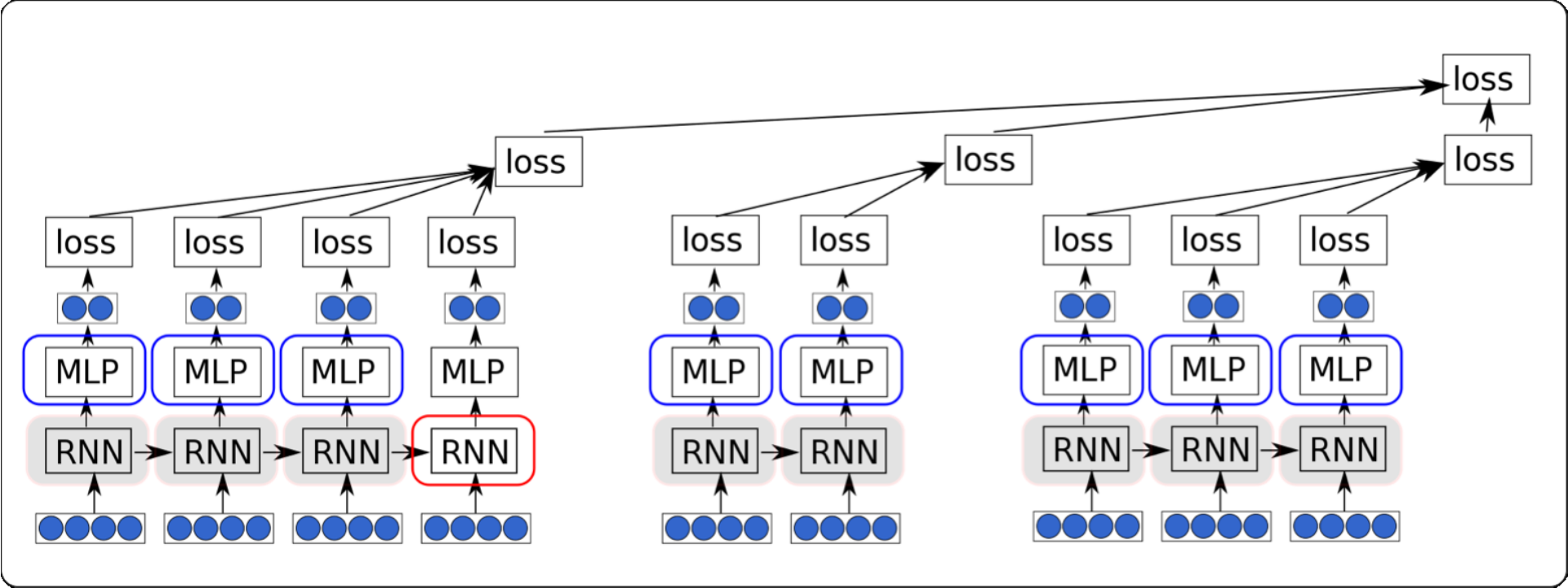


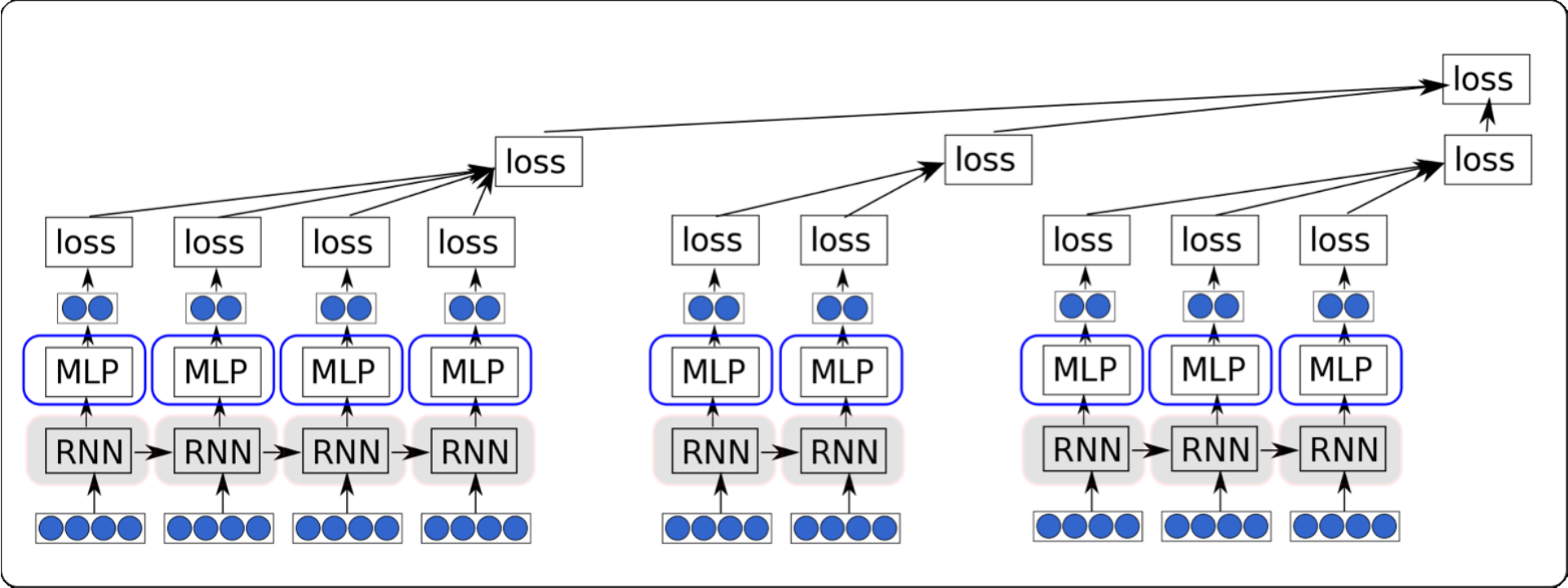


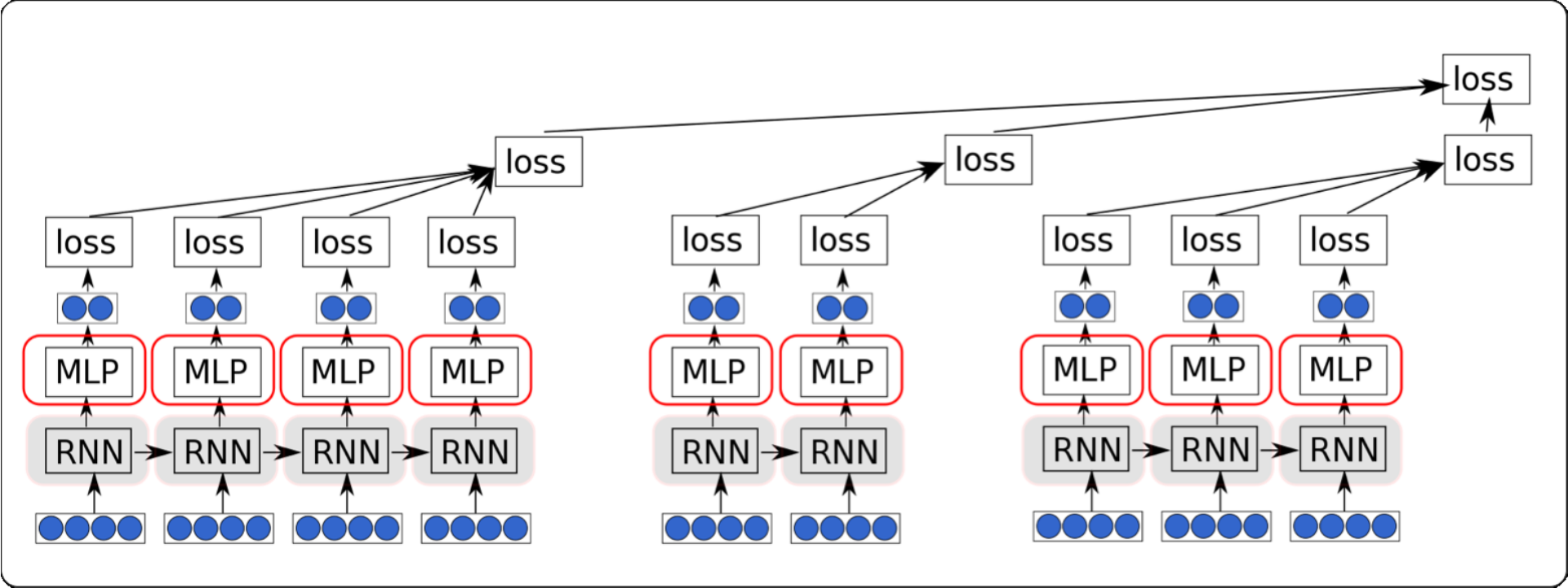


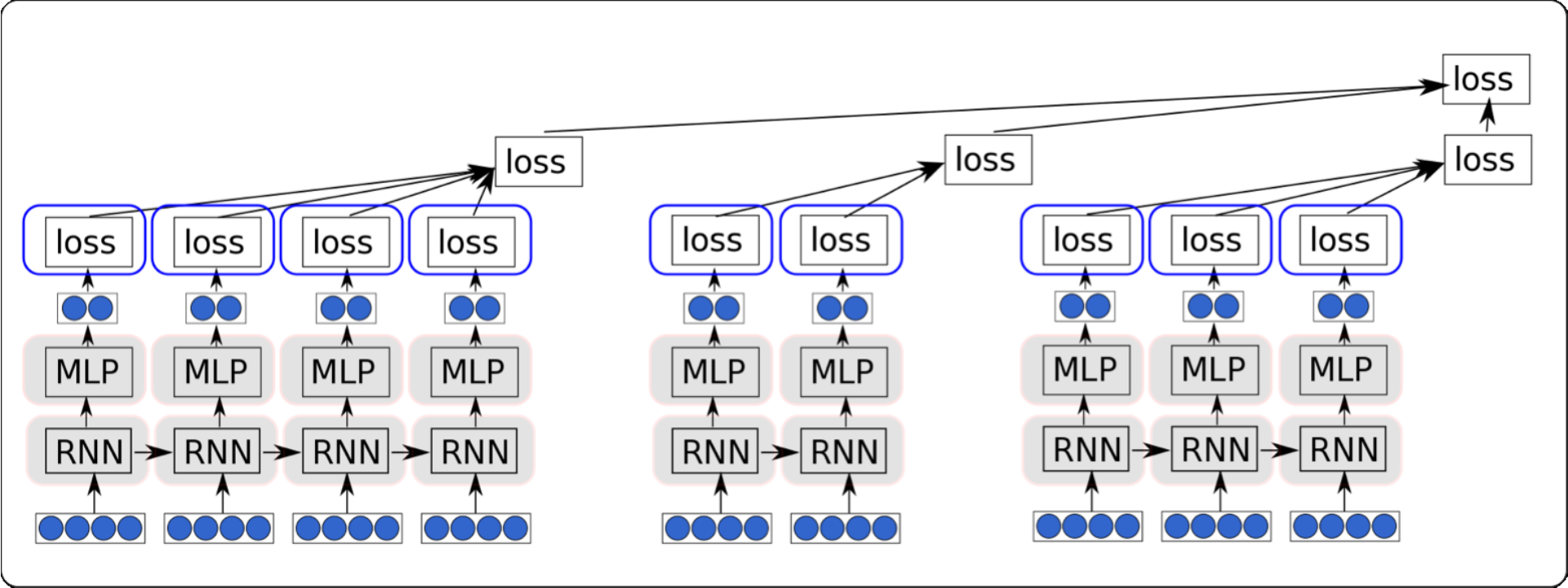


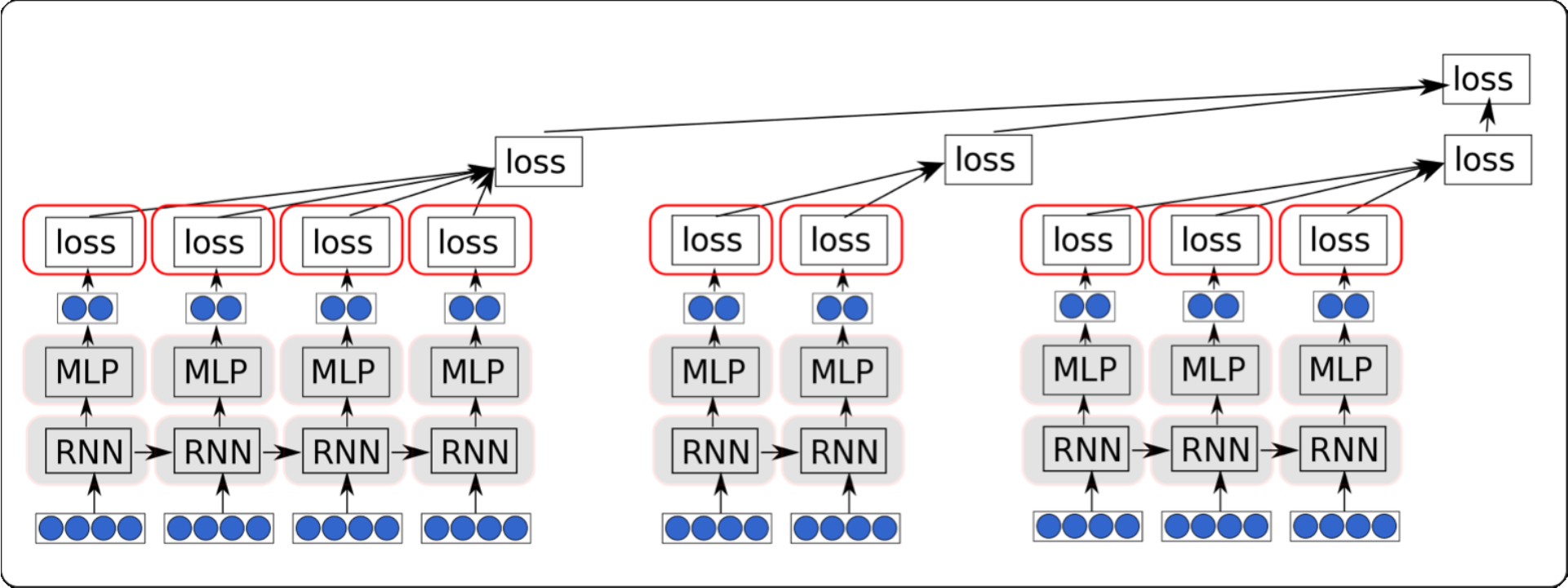


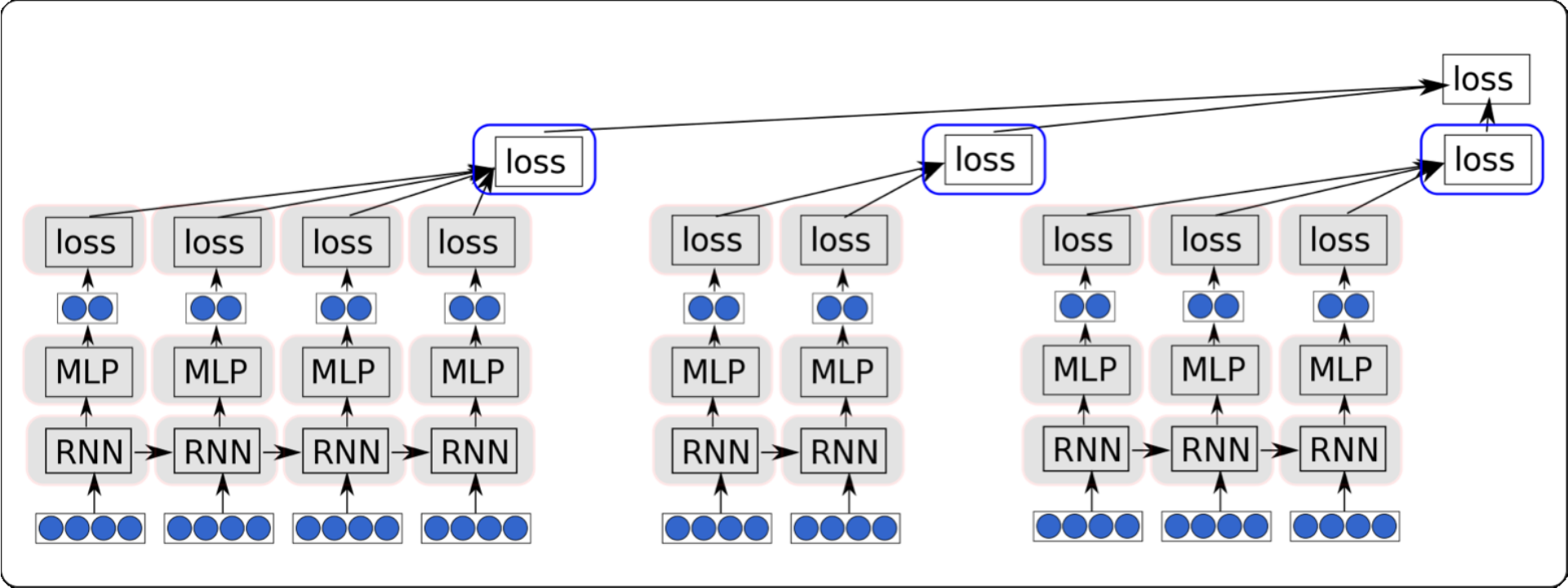


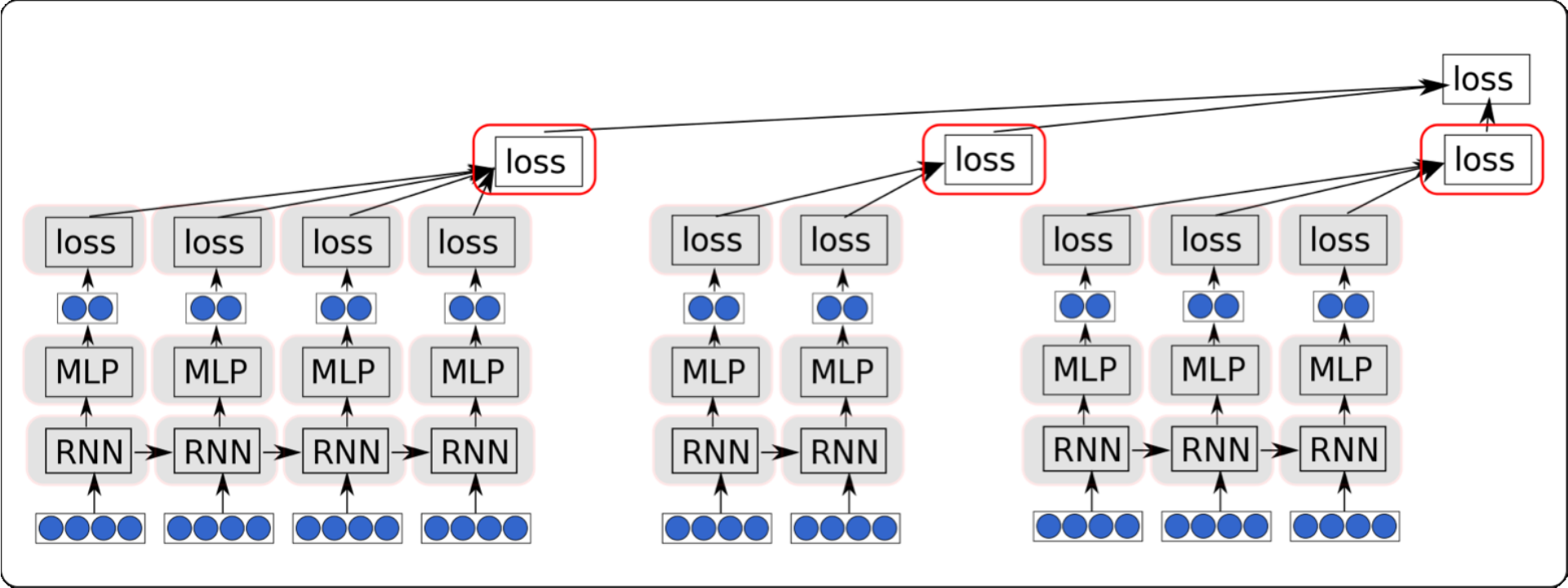


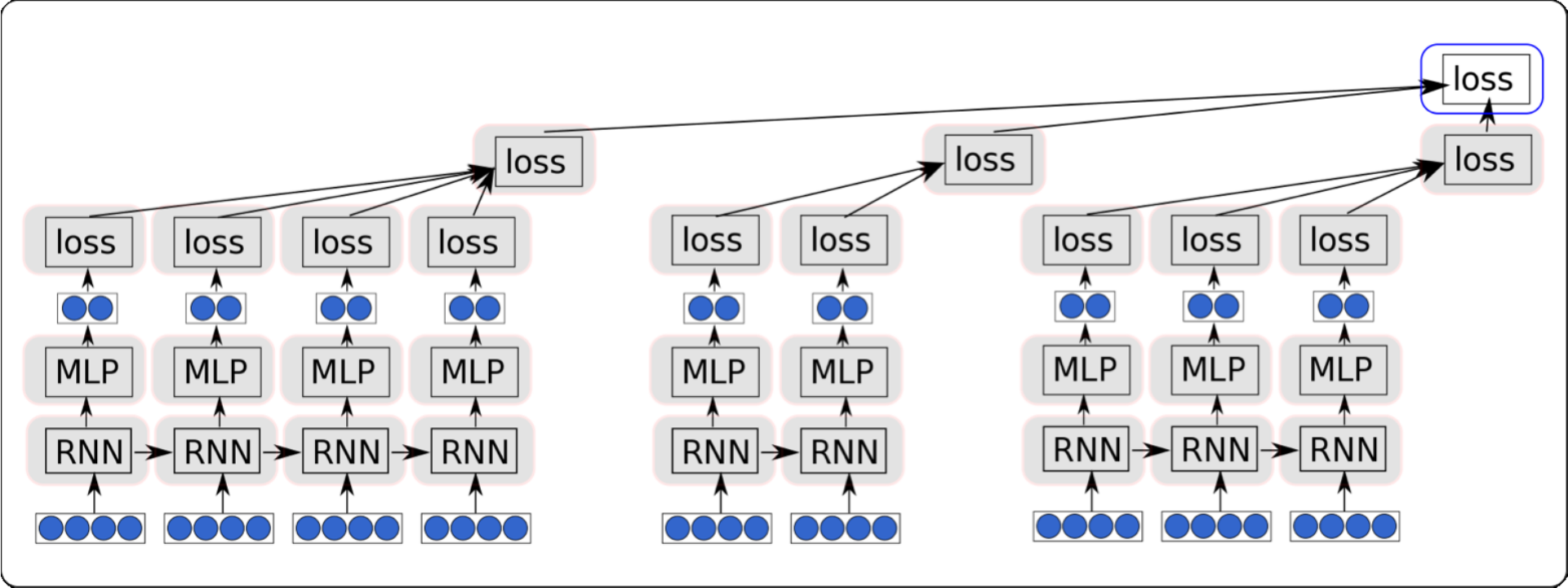


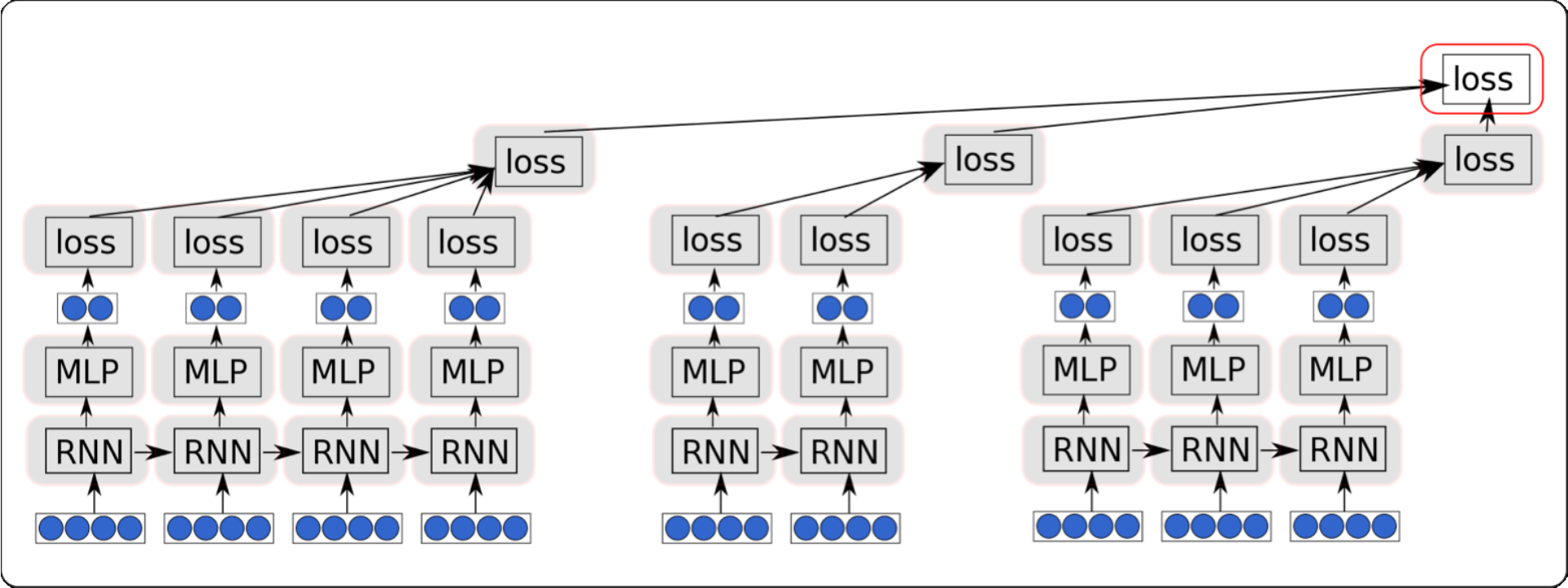


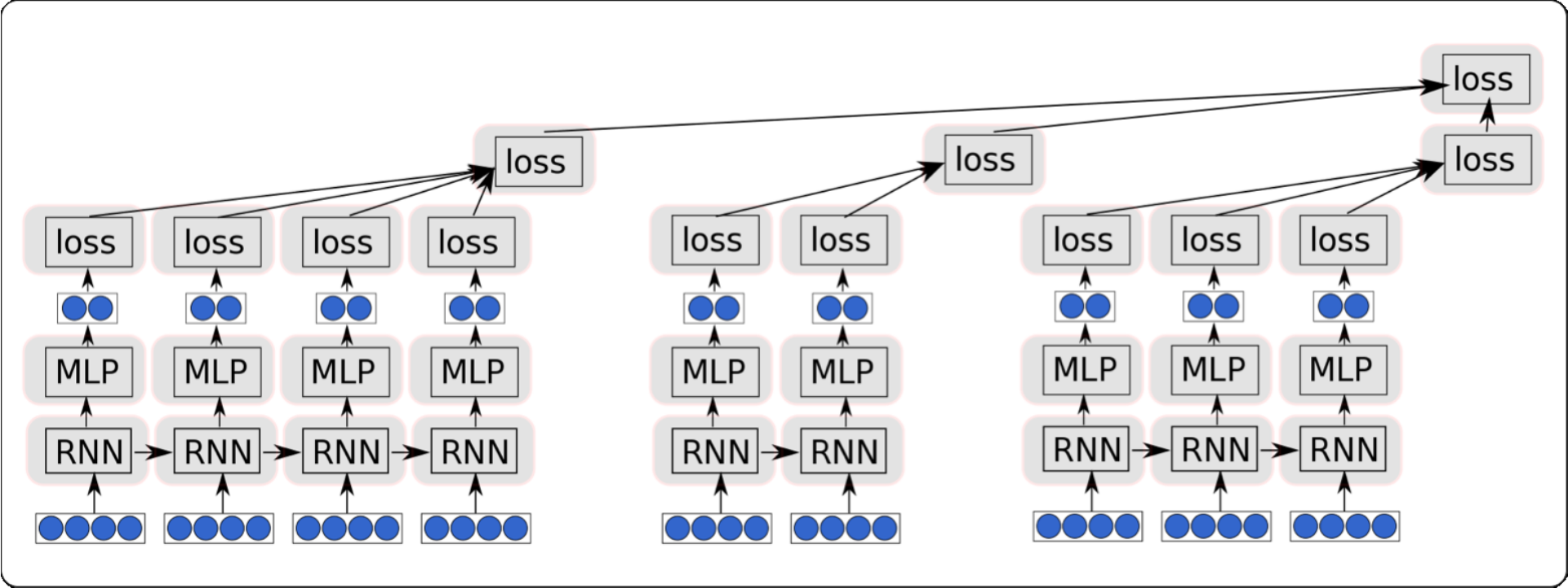


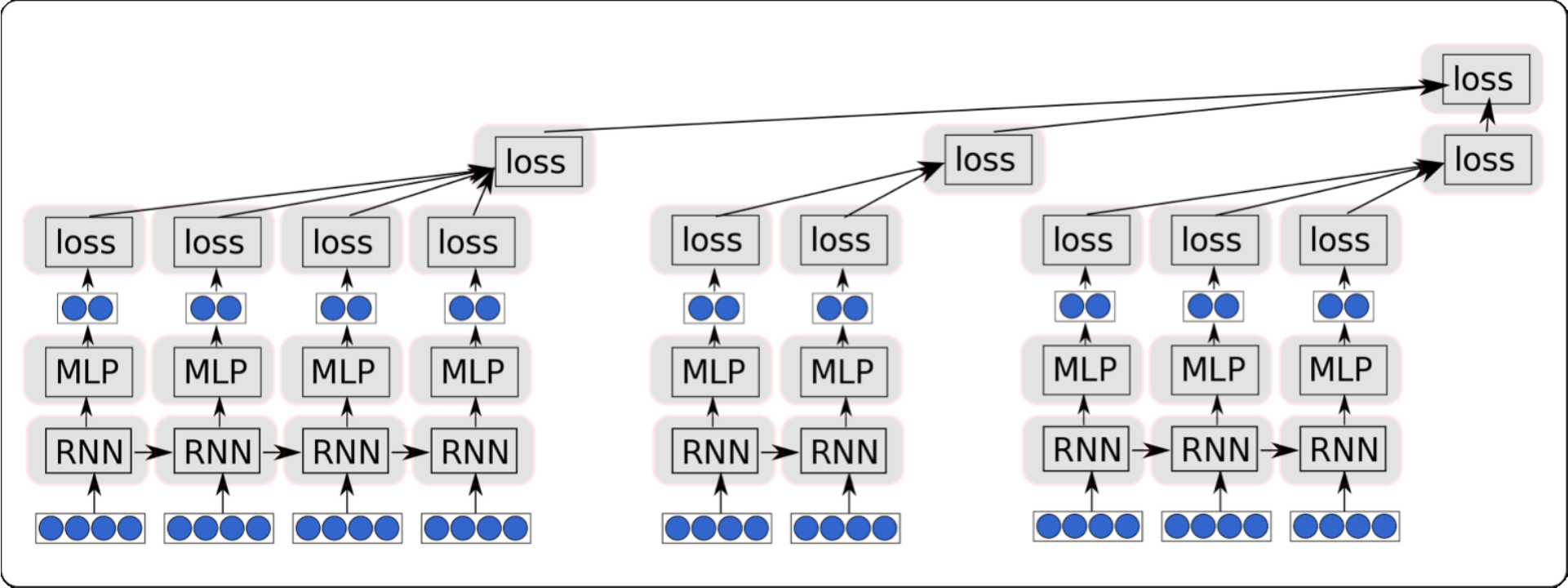












note: batching operations, not inputs.

Efficiency Considerations when Implementing an LSTM

$$R_{LSTM}(s_{j-1}, \mathbf{x}_j) = [\mathbf{c}_j; \mathbf{h}_j]$$

$$\mathbf{c}_j = \mathbf{c}_{j-1} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i}$$

$$\mathbf{h}_j = \tanh(\mathbf{c}_j) \odot \mathbf{o}$$

$$\mathbf{i} = \sigma(\mathbf{W}^{xi} \cdot \mathbf{x}_j + \mathbf{W}^{hi} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{f} = \sigma(\mathbf{W}^{xf} \cdot \mathbf{x}_j + \mathbf{W}^{hf} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{o} = \sigma(\mathbf{W}^{xo} \cdot \mathbf{x}_j + \mathbf{W}^{ho} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{g} = \tanh(\mathbf{W}^{xg} \cdot \mathbf{x}_j + \mathbf{W}^{hg} \cdot \mathbf{h}_{j-1})$$

Efficiency Considerations when Implementing an LSTM

$$R_{LSTM}(s_{j-1}, \mathbf{x}_j) = [\mathbf{c}_j; \mathbf{h}_j]$$

$$\mathbf{c}_j = \mathbf{c}_{j-1} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i}$$

$$\mathbf{h}_j = \tanh(\mathbf{c}_j) \odot \mathbf{o}$$

$$\mathbf{i} = \sigma(\mathbf{W}^{xi} \cdot \mathbf{x}_j + \mathbf{W}^{hi} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{f} = \sigma(\mathbf{W}^{xf} \cdot \mathbf{x}_j + \mathbf{W}^{hf} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{o} = \sigma(\mathbf{W}^{xo} \cdot \mathbf{x}_j + \mathbf{W}^{ho} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{g} = \tanh(\mathbf{W}^{xg} \cdot \mathbf{x}_j + \mathbf{W}^{hg} \cdot \mathbf{h}_{j-1})$$

all gates computations can be done in single mat-mat op.

Speed Trick 1: Don't Repeat Operations

- Something that you can do once at the beginning of the sentence, don't do it for every word!

Bad

```
for x in words_in_sentence:  
    vals.append( W * c + x )
```

Speed Trick 1:

Don't Repeat Operations

- Something that you can do once at the beginning of the sentence, don't do it for every word!

Bad

```
for x in words_in_sentence:  
    vals.append( W * c + x )
```

Good

```
W_c = W * c  
for x in words_in_sentence:  
    vals.append( W_c + x )
```

Speed Trick 2:

Reduce # of Operations

- e.g. can you combine multiple matrix-vector multiplies into a single matrix-matrix multiply? Do so!

Bad

```
for x in words_in_sentence:  
    vals.append( W * x )  
val = dy.concatenate(vals)
```

Speed Trick 2:

Reduce # of Operations

- e.g. can you combine multiple matrix-vector multiplies into a single matrix-matrix multiply? Do so!

Bad

```
for x in words_in_sentence:  
    vals.append( W * x )  
val = dy.concatenate(vals)
```

Good

```
X = dy.concatenate_cols(words_in_sentence)  
val = W * X
```

Speed Trick 3:

Reduce CPU-GPU Data Movement

- Try to **avoid memory moves** between CPU and GPU.
- When you do move memory, try to do it as early as possible (GPU operations are asynchronous)

Bad

```
for x in words_in_sentence:  
    # input data for x  
    # do processing
```

Speed Trick 3:

Reduce CPU-GPU Data Movement

- Try to **avoid memory moves** between CPU and GPU.
- When you do move memory, try to do it as early as possible (GPU operations are asynchronous)

Bad

```
for x in words_in_sentence:  
    # input data for x  
    # do processing
```

Good

```
# input data for whole sentence  
for x in words_in_sentence:  
    # do processing
```