

DeepBind

Predicting the sequence specificities of DNA- and RNA-binding proteins by deep learning

Nature Biotechnology 33, 831–838 (2015) doi:10.1038/nbt.3300

Method Details

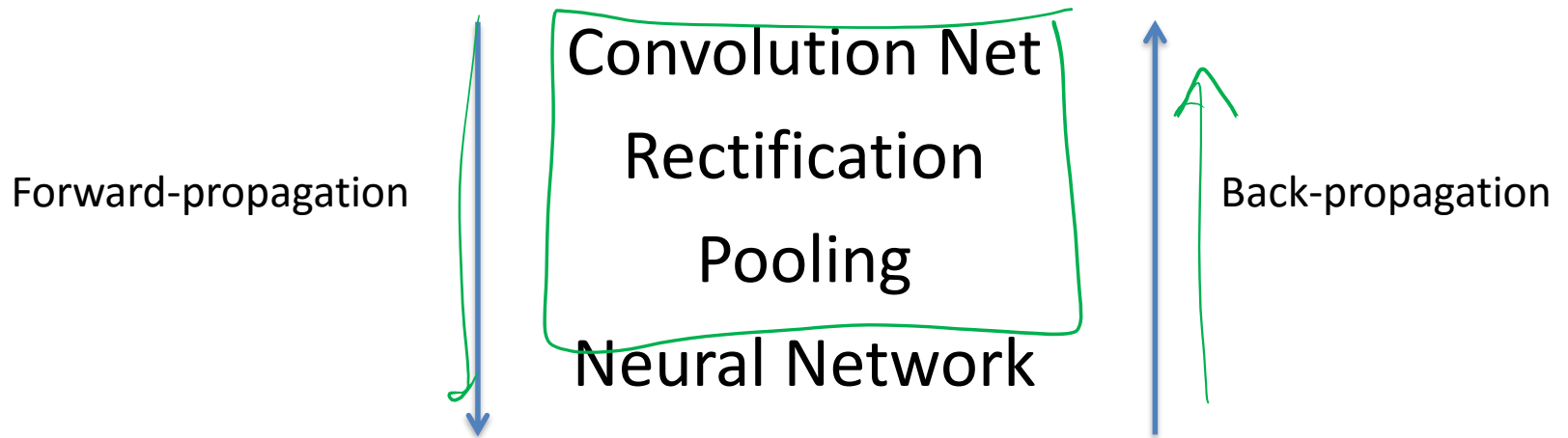
Hayan Lee

Research Fellow @ Simons Institute for the Theory of Computing, UC Berkeley

Simons Postdoctoral Fellow @ DOE JGI, Lawrence Berkeley National Laboratory

Outline

- Overview
- Forward propagation
- Backward propagation



Overview



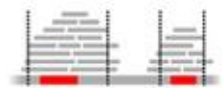
1. High-throughput experiments



PBM



SELEX

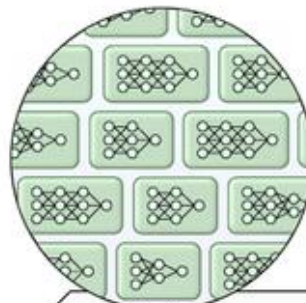


CHIP/CLIP

Large-scale
data sets

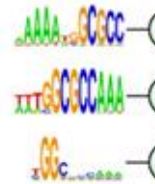
2. Massively parallel deep learning

Automatic model training

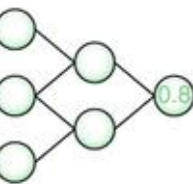


GPU server

New
motifs



Prediction
network



DeepBind
models

3. Community needs



Gene regulation



Precision medicine

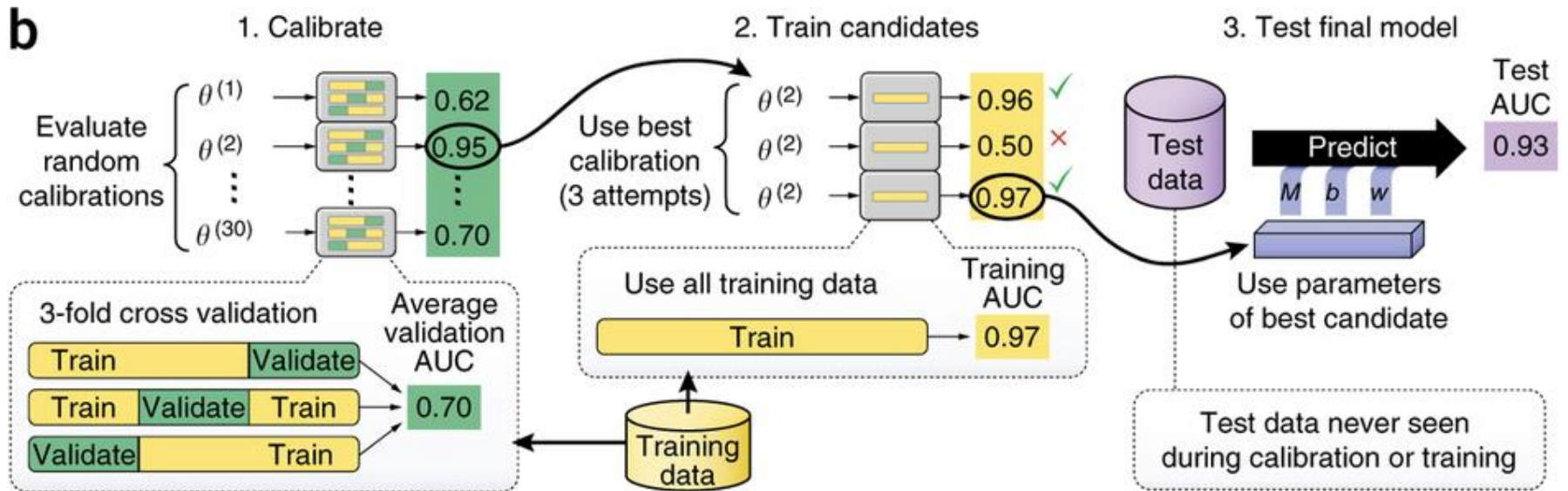
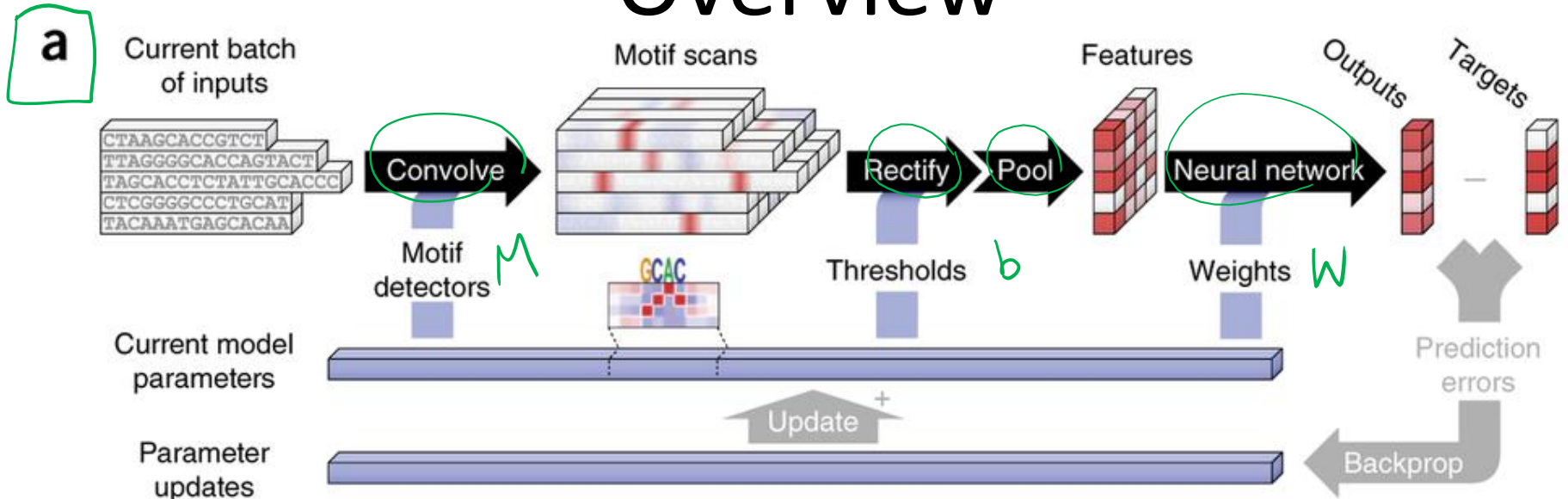


Detect binding sites

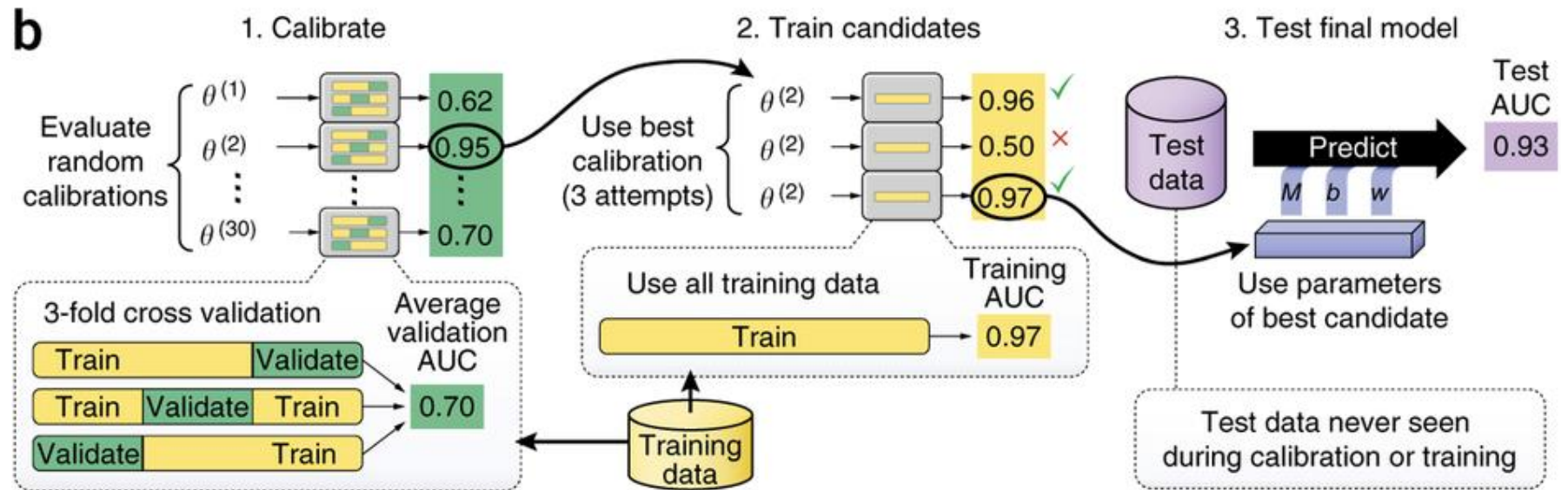
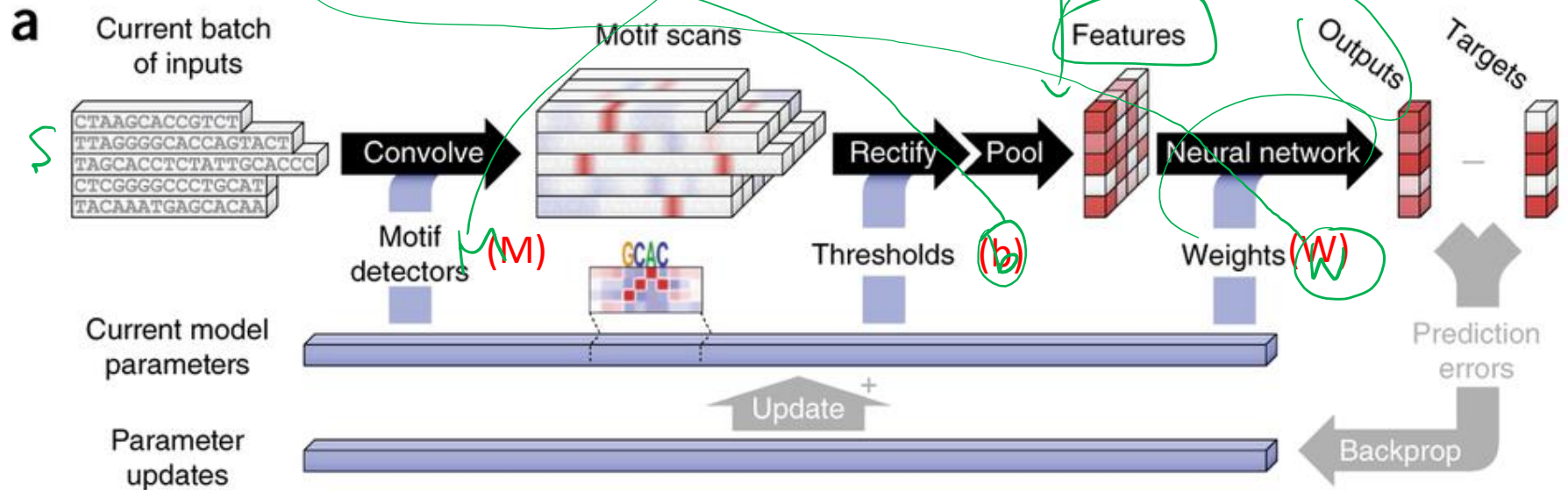
14–101 nt

PBM : Protein binding microarrays

Overview



$$f(s) = \text{net}_w(\text{pool}(\text{rect}_b(\text{conv}_M(s)))) \quad (M, b \text{ and } w \text{ are trainable parameters})$$



Convolution

$$f(s) = net_w(pool(rect_b(conv_M(s))))$$

$s = (s_1, \dots, s_n)$ Sequences can have varying lengths (14–101 nt in our experiments),

$conv_M(s)$ The convolution stage scans a set of motif detectors with parameters M across the sequence. Motif detector M_k is a $m \times 4$ matrix, much like a PWM (Position Weight Matrices) of length m but without requiring coefficients to be probabilities or log odds ratios.

Convolution : Encoding & Padding

Input Sequence: ATGG

$$S_{i,j} = \begin{cases} .25 & \text{if } s_{i-m+1} = N \text{ or } i < m \text{ or } i > n - m \\ 1 & \text{if } s_{i-m+1} = j^{\text{th}} \text{ base in (A, C, G, T)} \\ 0 & \text{otherwise} \end{cases}$$

n is sequence length,
 m is motif length

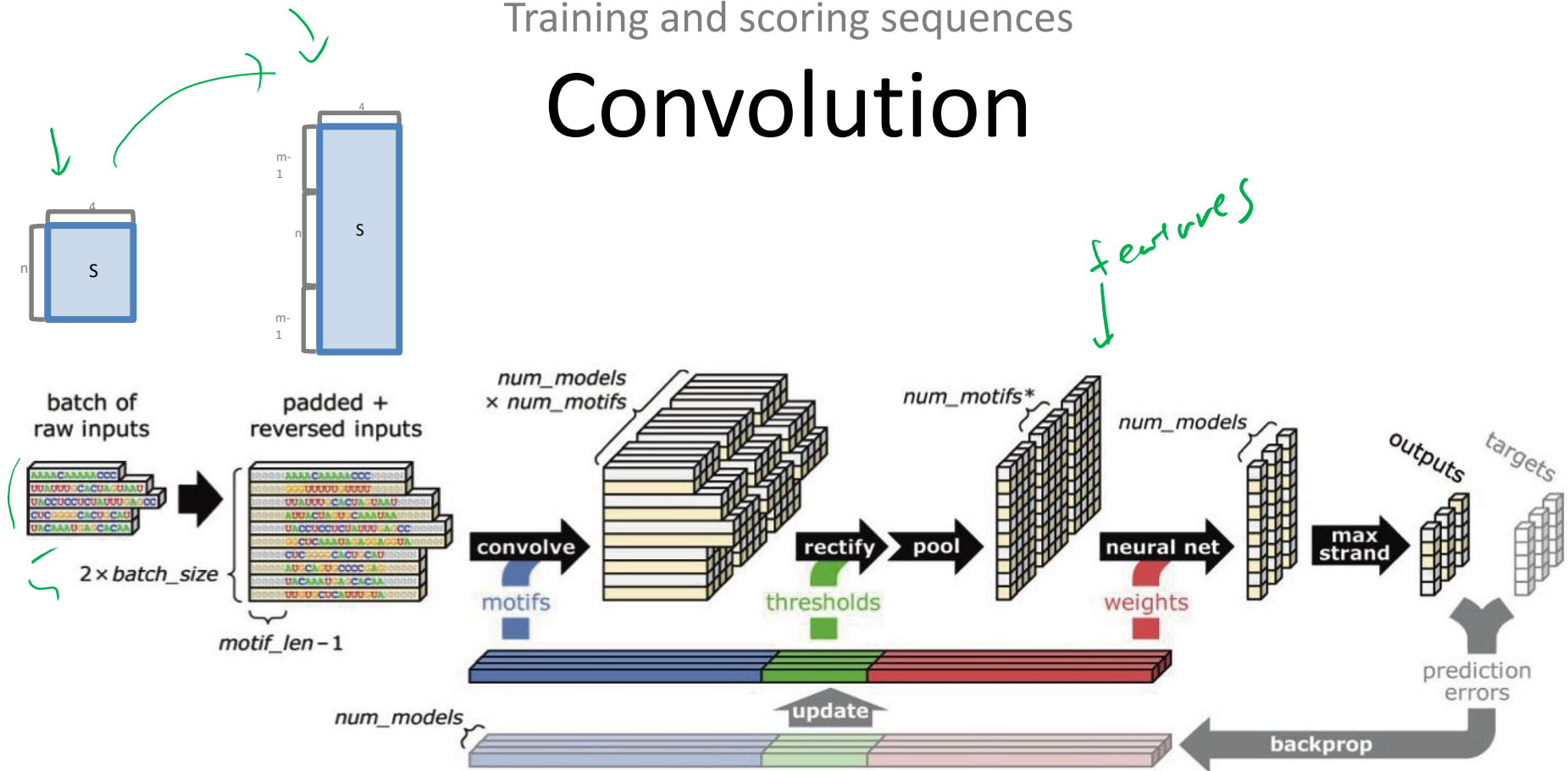
↓

$$S = \begin{matrix} & \begin{matrix} A & C & G & T \end{matrix} \\ \begin{matrix} m-1 \\ n \\ n-1 \end{matrix} \left[\begin{array}{cccc} .25 & .25 & .25 & .25 \\ .25 & .25 & .25 & .25 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ .25 & .25 & .25 & .25 \\ .25 & .25 & .25 & .25 \end{array} \right] & \begin{matrix} \Sigma = 1 \\ n-1 \\ A \\ T \\ G \\ G \\ n-1 \end{matrix} \end{matrix}$$

$$\text{Size}(S) = n + 2(m-1)$$

Training and scoring sequences

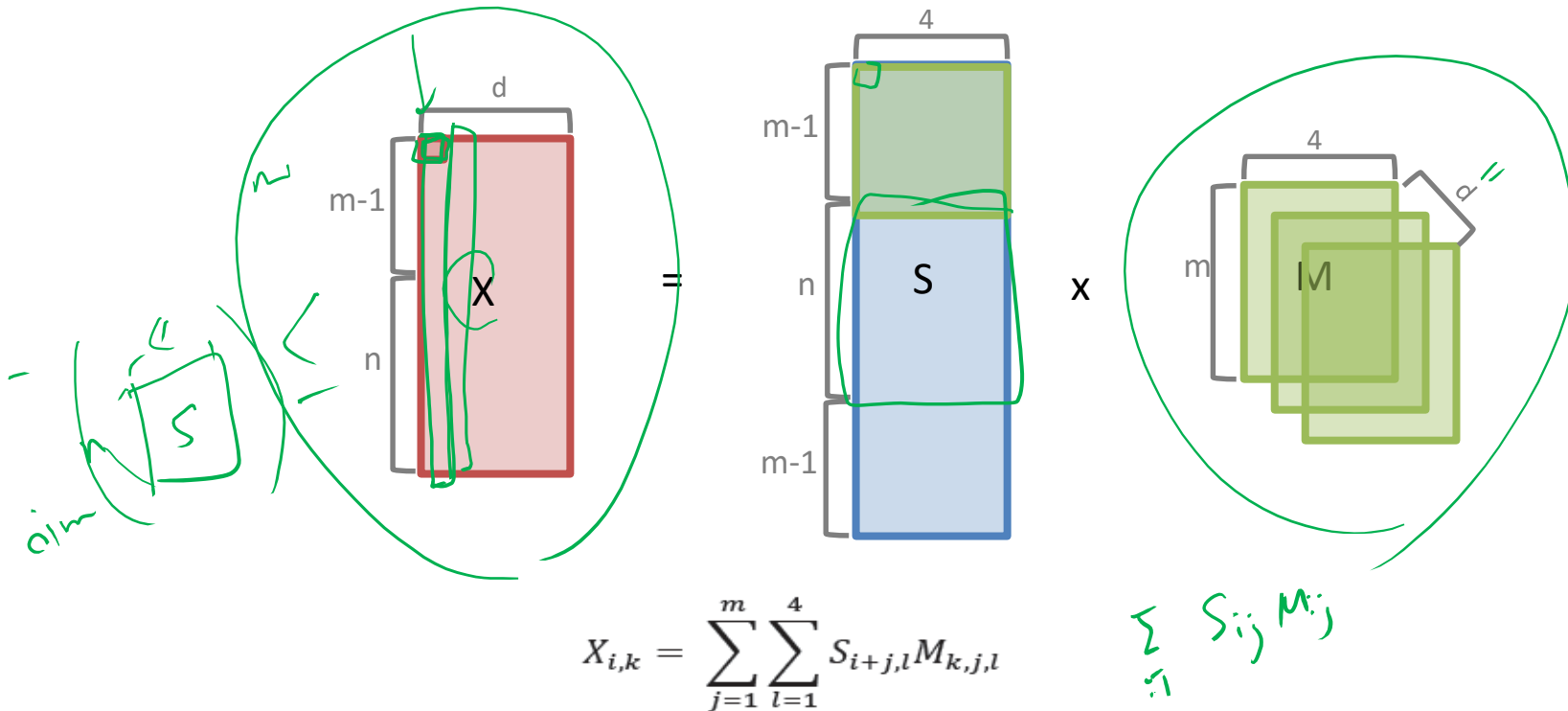
Convolution



hyper parameters

batch_size=5, motif_len=6, num_motifs=4, num_models=3

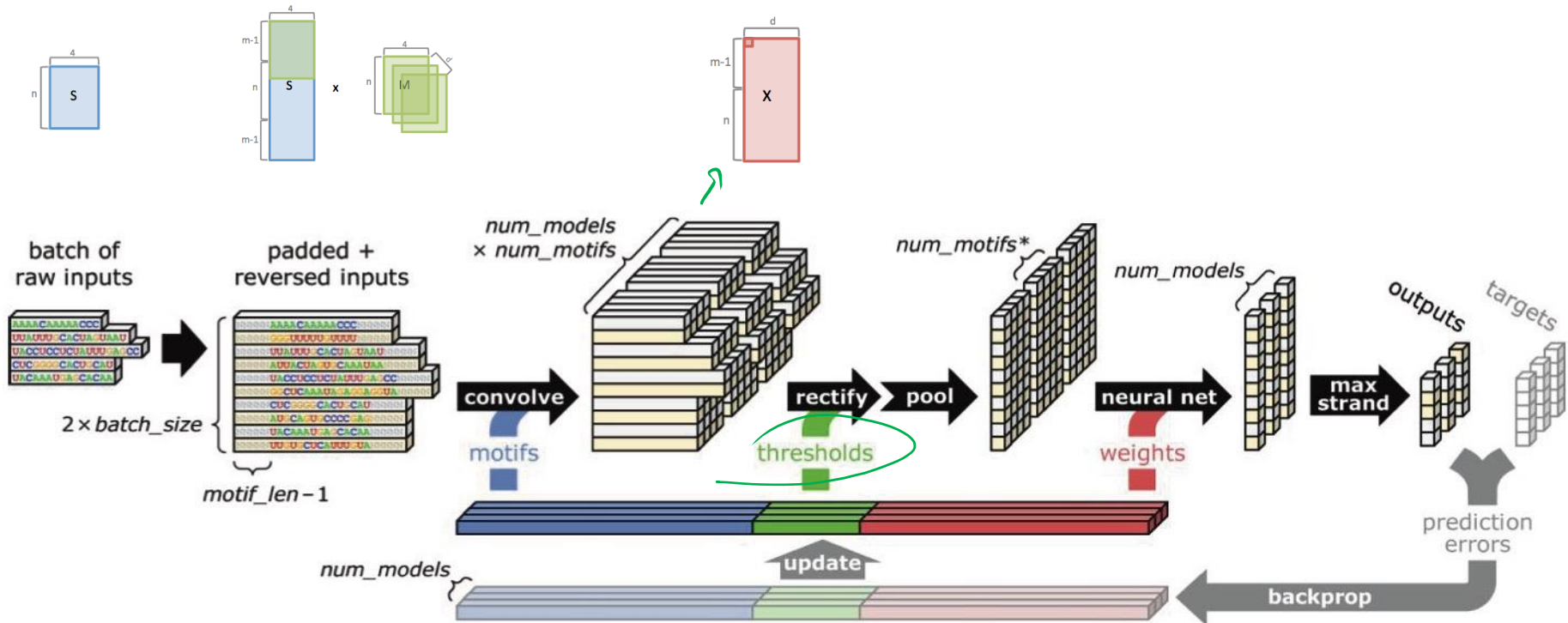
Convolution



- For all $1 \leq i \leq n+m-1$ and $1 \leq k \leq d$, where d is the number of motif detectors
- S (input sequences) : $(n+2m-2) \times 4$
- M (motif detectors) : $m \times 4 \times d$
- X (output) : $(n+m-1) \times d$

Training and scoring sequences

Rectification



batch_size=5, motif_len=6, num_motifs=4, num_models=3

Rectification

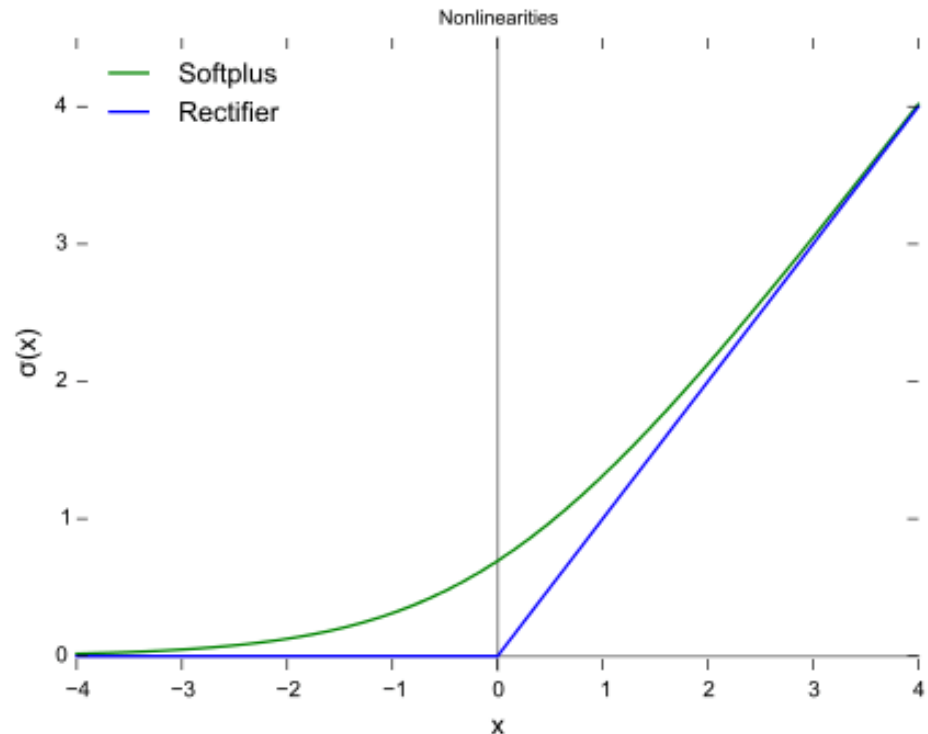
$$f(s) = net_w(pool(rect_b(conv_M(s))))$$

$rect_b$

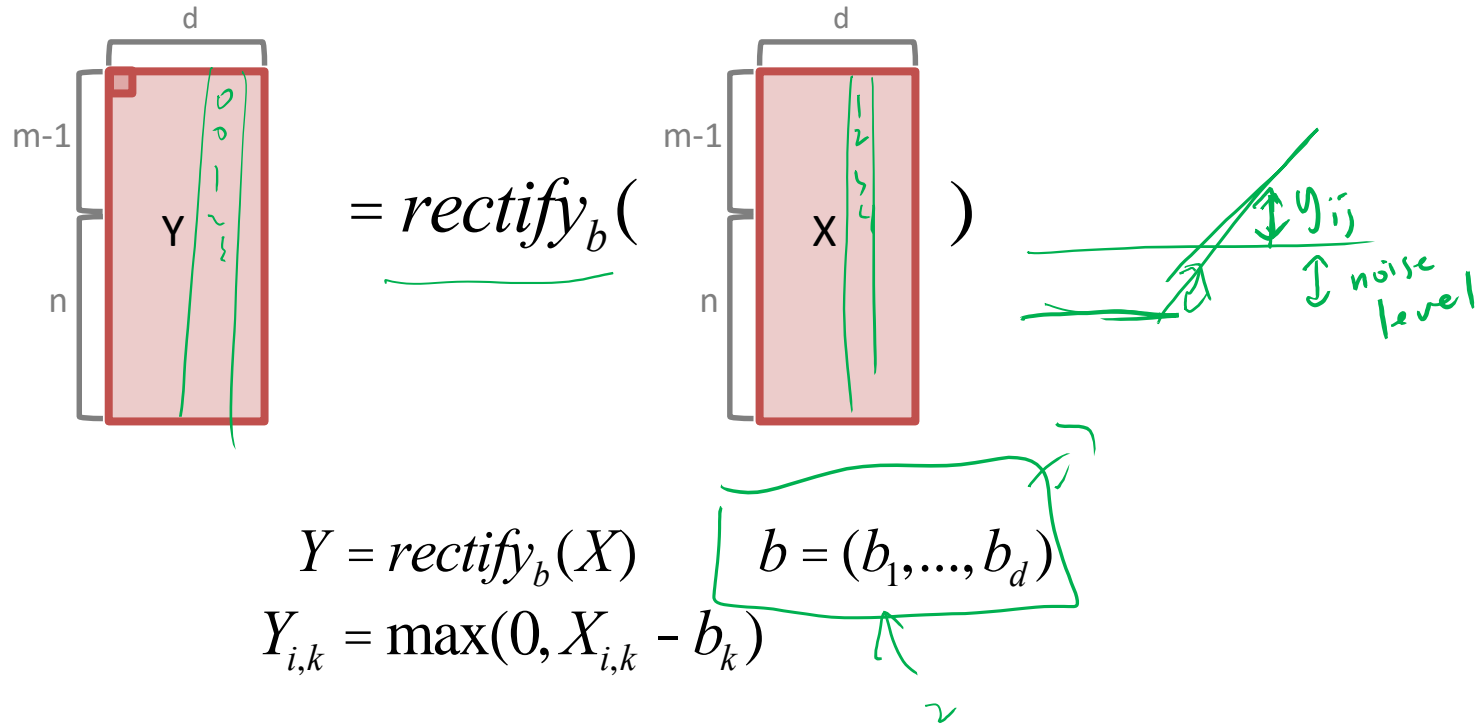
The rectification stage isolates positions with a good pattern match by shifting the response of detector M_k by b_k and clamping all negative values to zero.

→ $f(x) = \ln(1 + e^x)$

→ $f(x) = \max(0, x)$



Rectification



- For all $1 \leq i \leq n+m-1$ and $1 \leq k \leq d$, where d is the number of motif detectors
- X (input) : $(n+m-1) \times 4$
- Y (output) : $(n+m-1) \times 4$

Pooling



feature
selection

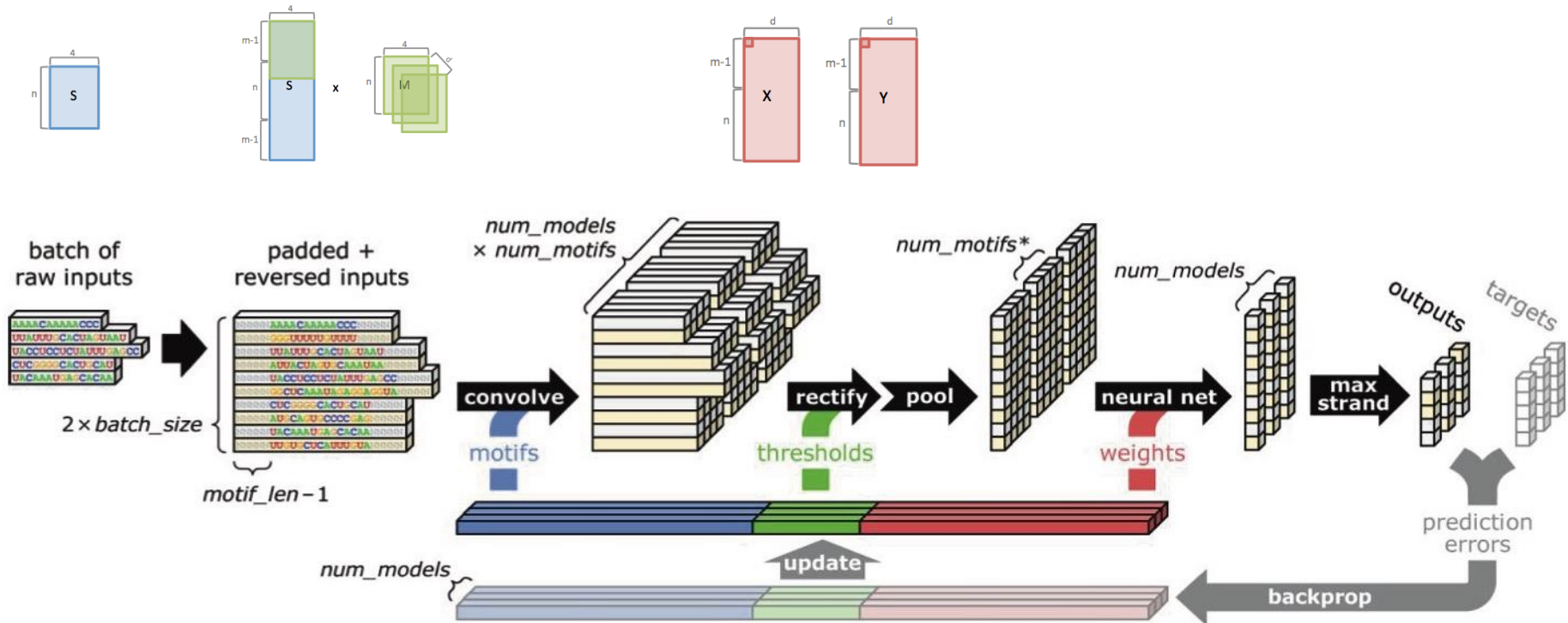
$$f(s) = net_w(pool(rect_b(conv_M(s))))$$

pool

The pooling stage computes the maximum and average of each motif detector's rectified response across the sequence; maximizing helps to identify the presence of longer motifs, whereas averaging helps to identify cumulative effects of short motifs, and the contribution of each is determined automatically by learning.

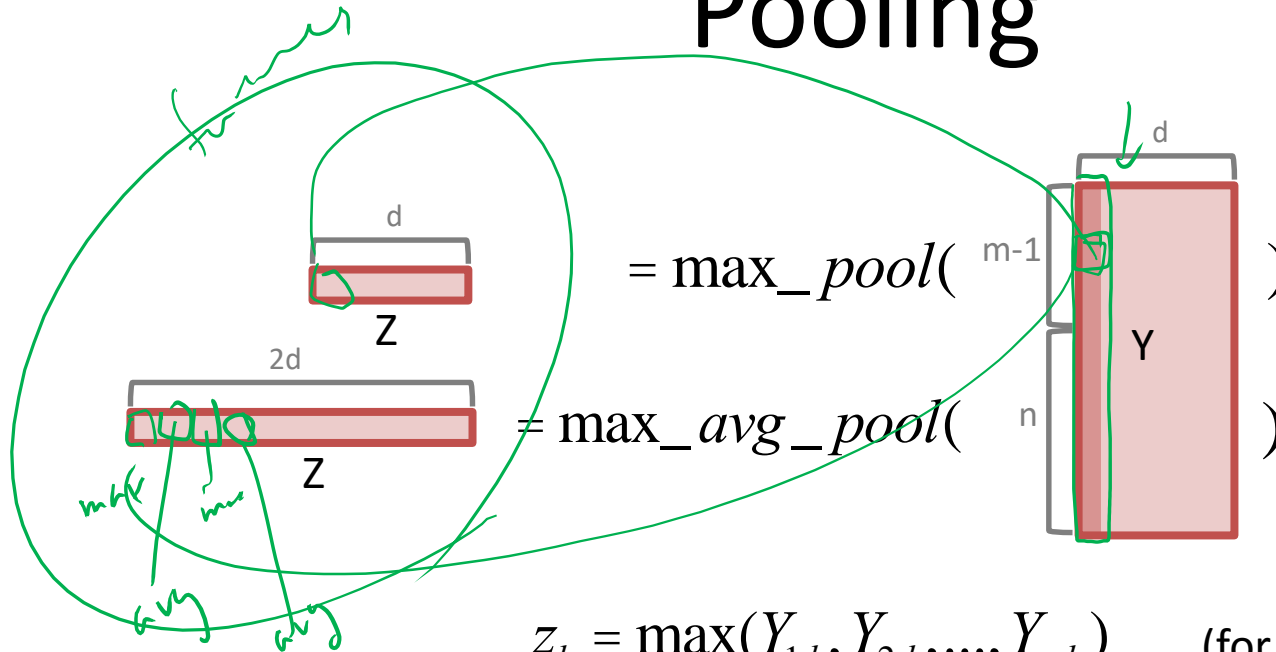
Training and scoring sequences

Pooling



batch_size=5, motif_len=6, num_motifs=4, num_models=3

Pooling



$$= \max_pool(\begin{matrix} m-1 \\ Y \end{matrix})$$

$$= \max_avg_pool(\begin{matrix} n \\ Y \end{matrix})$$

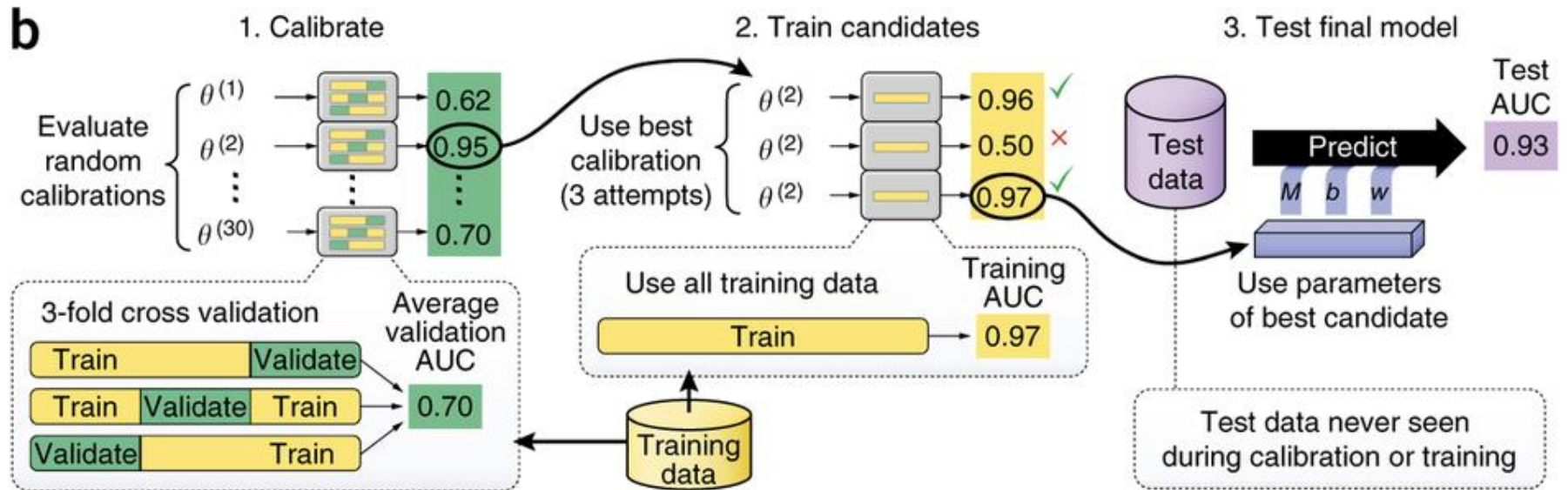
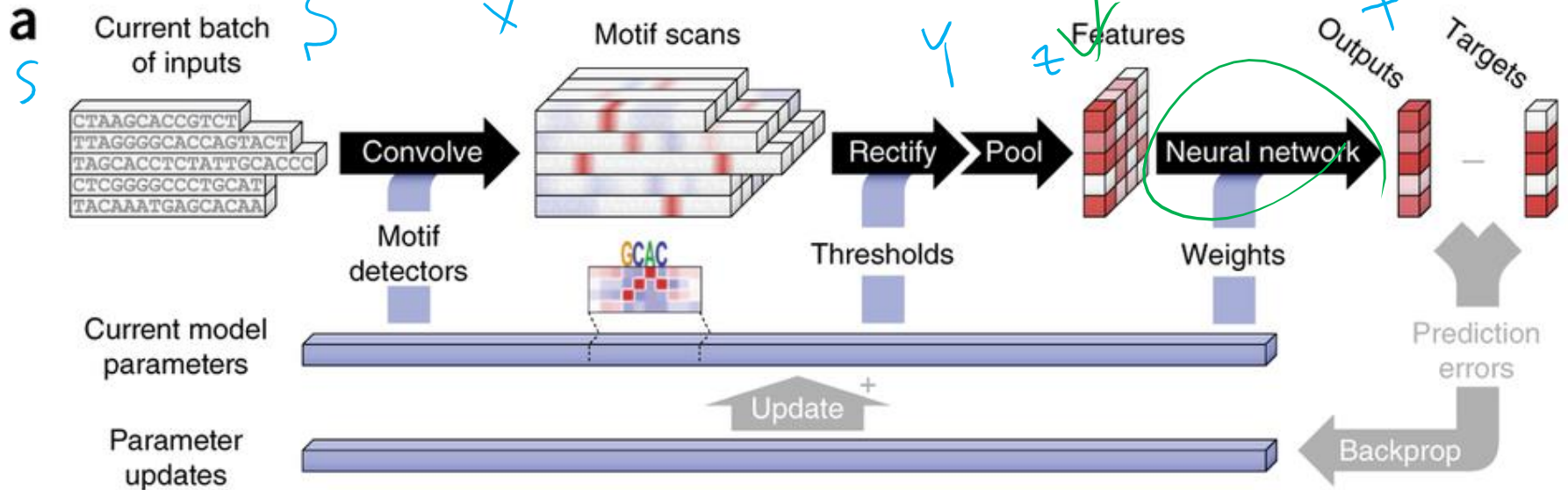
$$z_k = \max(Y_{1,k}, Y_{2,k}, \dots, Y_{n,k}) \quad (\text{for DNA binding})$$

$$z_{2k+0} = \max(Y_{1,k}, Y_{2,k}, \dots, Y_{n,k}) \quad (\text{for RNA binding})$$

$$z_{2k+1} = \text{avg}(Y_{1,k}, Y_{2,k}, \dots, Y_{n,k})$$

- $1 \leq k \leq d$, where d is the number of motif detectors
- Y (input) : $(n+m-1) \times 4$
- Z (output) : $1 \times d$ (for DNA binding)
- Z (output) : $1 \times 2d$ (for RNA binding, RBP model)

Overview



Training and scoring sequences

Neural network

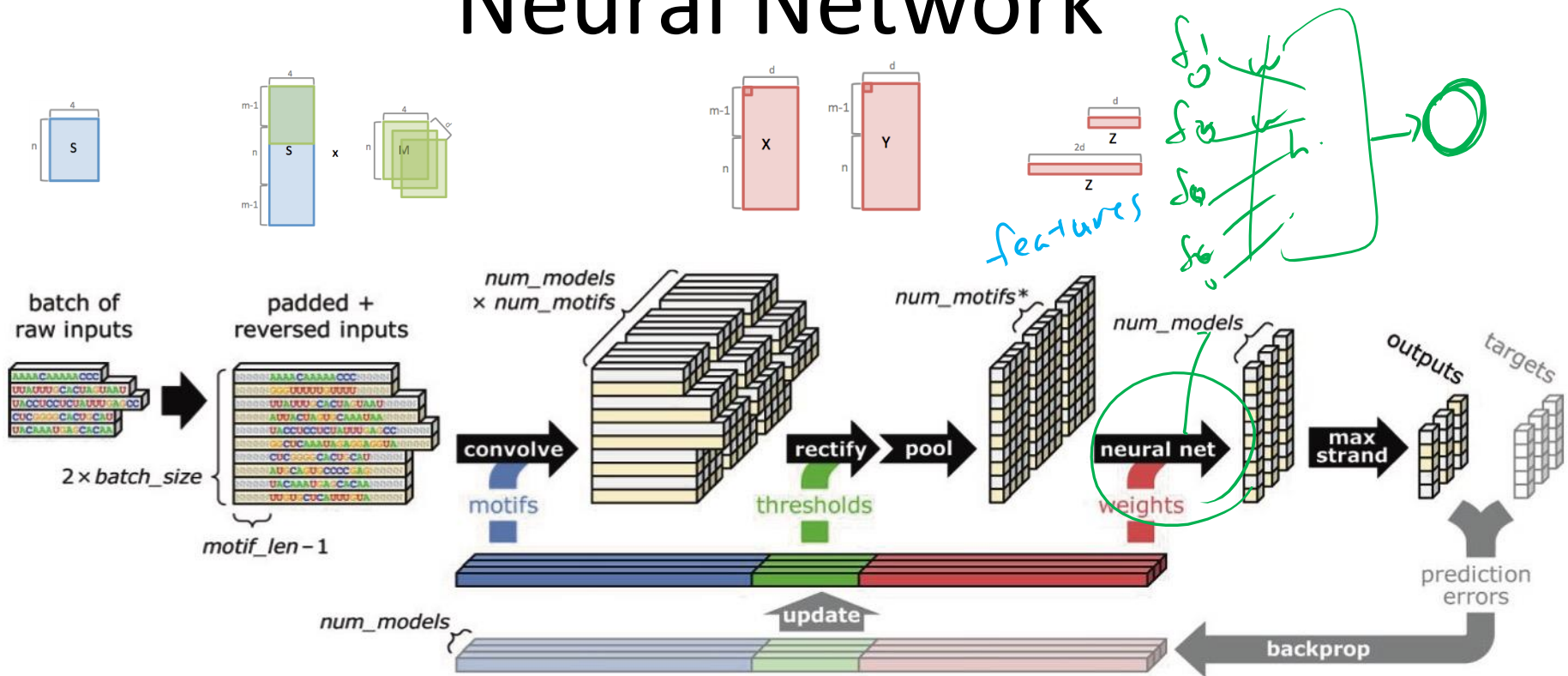
$$f(s) = \textit{net}_w(\textit{pool}(\textit{rect}_b(\textit{conv}_M(s))))$$

\textit{net}_w

Output vector z from pooling plays a role of feature vector.
These values are fed into a nonlinear neural network with weights \mathbf{W} ,
which combines the responses to produce a score

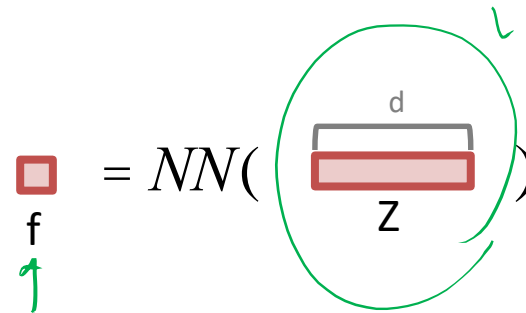
Training and scoring sequences

Neural Network



batch_size=5, motif_len=6, num_motifs=4, num_models=3

Neural Network



Without hidden layer

$$p = w_{d+1} + \sum_{k=1}^d w_k z_k$$

One hidden layer with 32 rectified-linear units(ReLU)

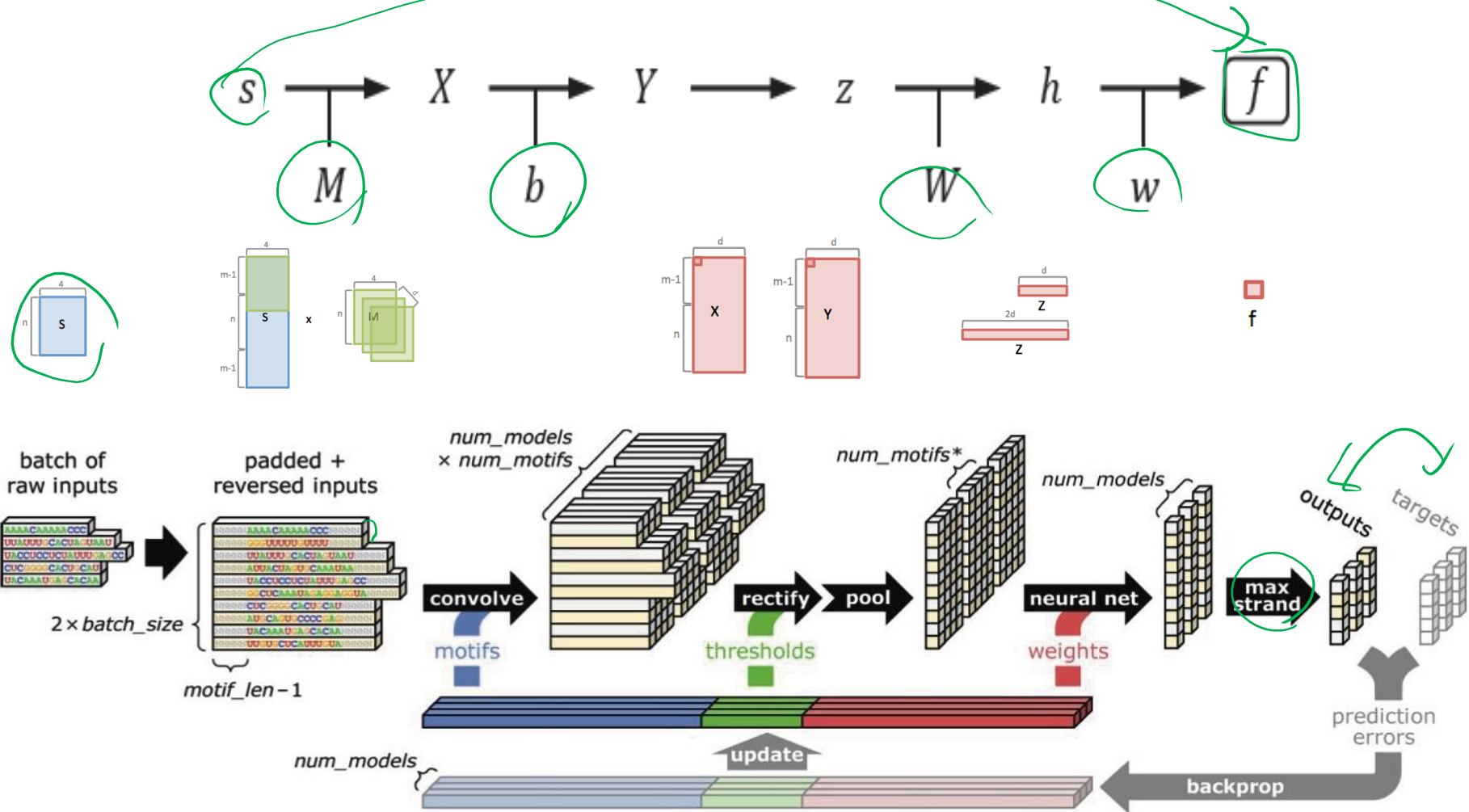
$$h_j = \max(0, w_{j,d+1} + \sum_{k=1}^d w_{j,k} z_k)$$

$$f = w_{33} + \sum_{j=1}^{32} w_j z_j$$

- $1 \leq k \leq d$, where d is the number of motif detectors
- Z (input) : $1 \times d$ (for DNA binding)
- Z (input) : $1 \times 2d$ (for RNA binding, RBP model)
- P or f : a scalar output score

Training and scoring sequences

Forward propagation

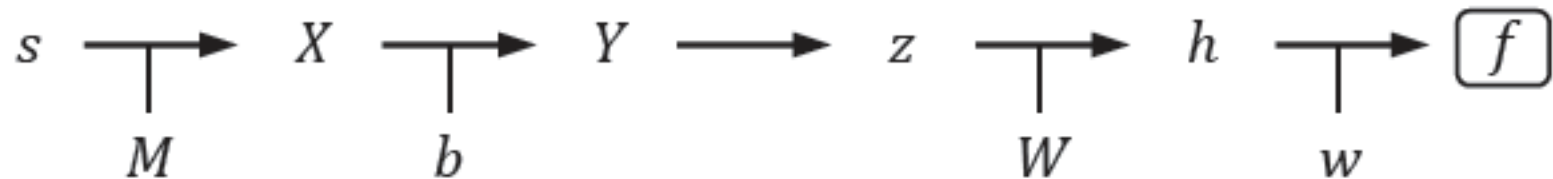


batch_size=5, motif_len=6, num_motifs=4, num_models=3

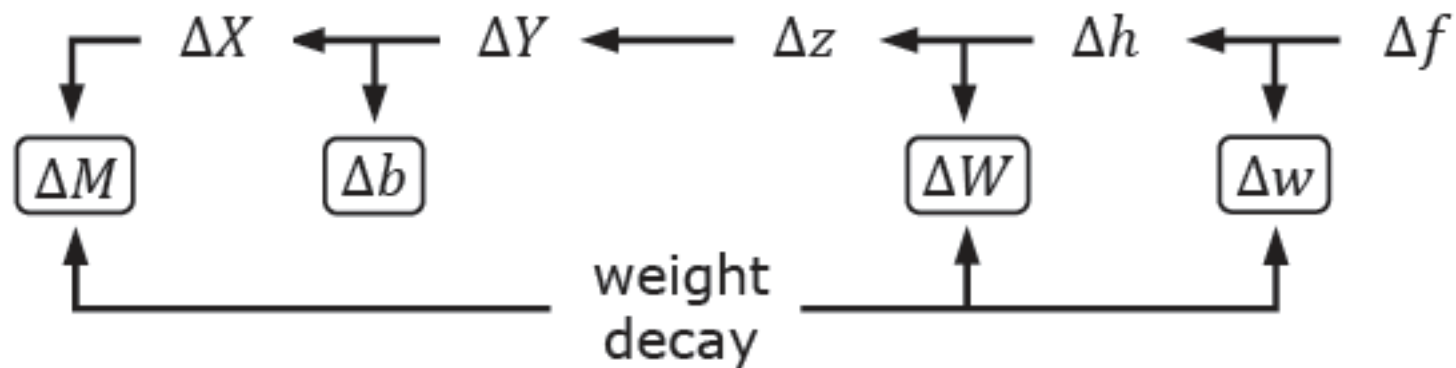
Training and scoring sequences

Back-propagation

Forward Propagation

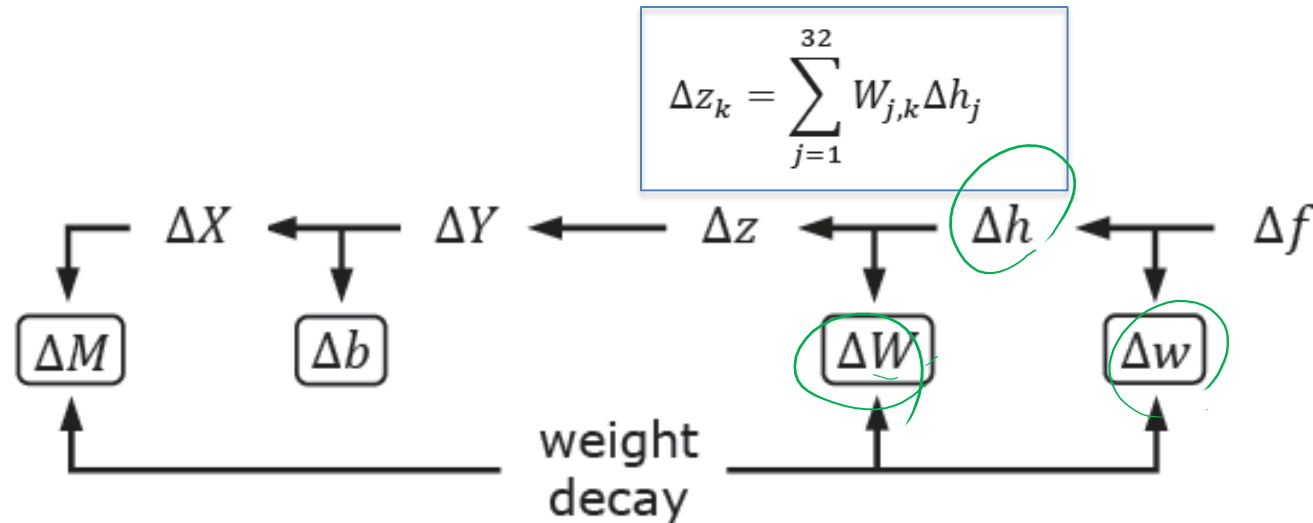


Back-propagation



Contributes to the gradients but independent of $\mathcal{D}f$

Back-propagation : NN



$$\begin{aligned} \Delta W_{j,k} &= z_k \Delta h_j \quad \text{for } k = 1 \dots d \\ \Delta W_{j,d+1} &= \Delta h_j \end{aligned}$$

Neural network with drop out

$$\begin{aligned} \Delta w_j &= m_j h_j \Delta f \quad \text{for } j = 1 \dots 32 \\ \Delta w_{33} &= \Delta f \end{aligned}$$

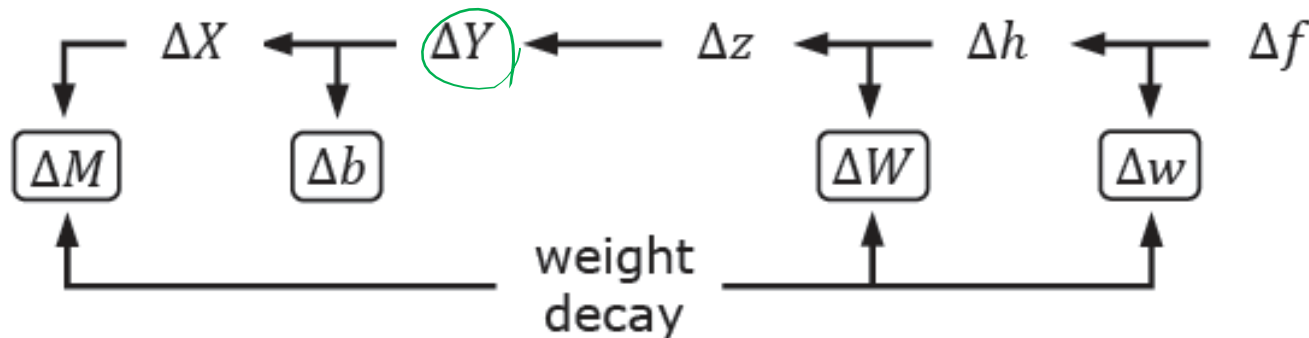
$m_j = \{0,1\}$, mask value for dropout

$$\Delta h_j = \begin{cases} m_j w_j \Delta f & \text{if } h_j > 0 \\ 0 & \text{otherwise} \end{cases}$$

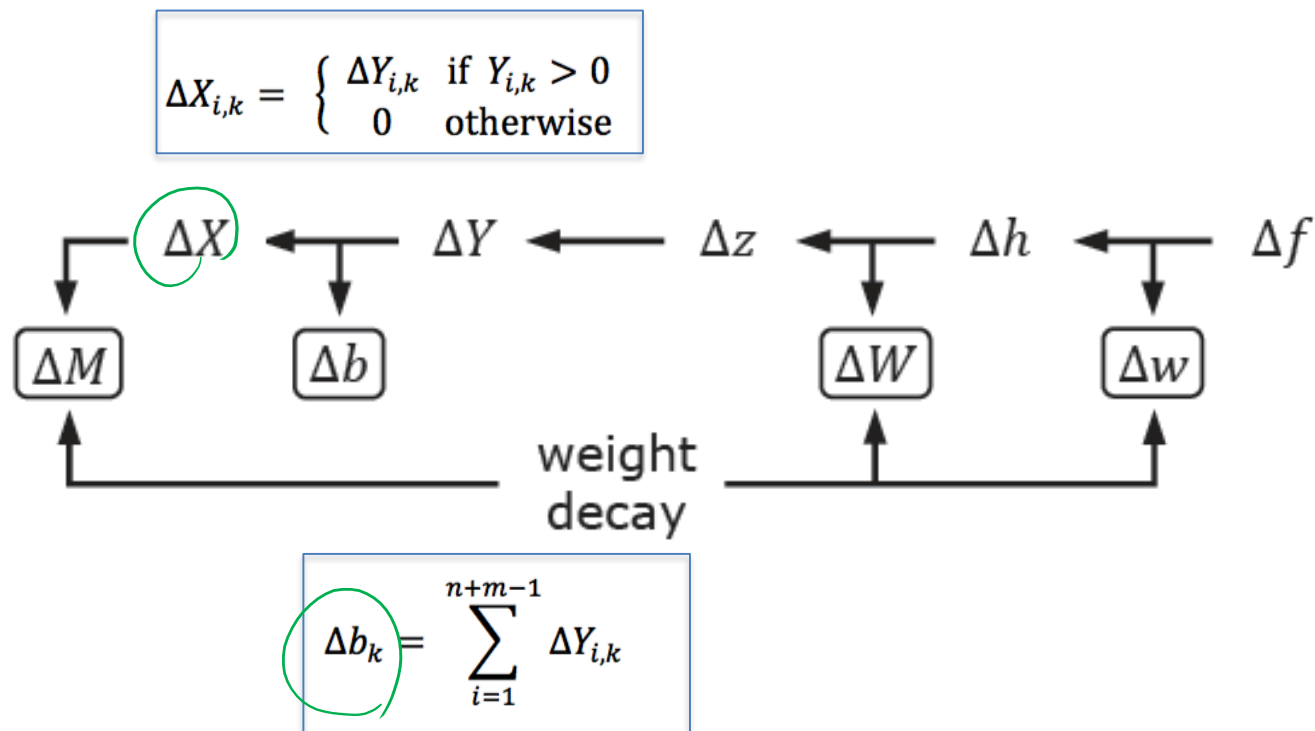
Back-propagation : Pooling

$$\Delta Y_{i,k} = \begin{cases} \Delta z_k & \text{if } i = \operatorname{argmax}(Y_{1,k}, \dots, Y_{n,k}) \\ 0 & \text{otherwise.} \end{cases}$$

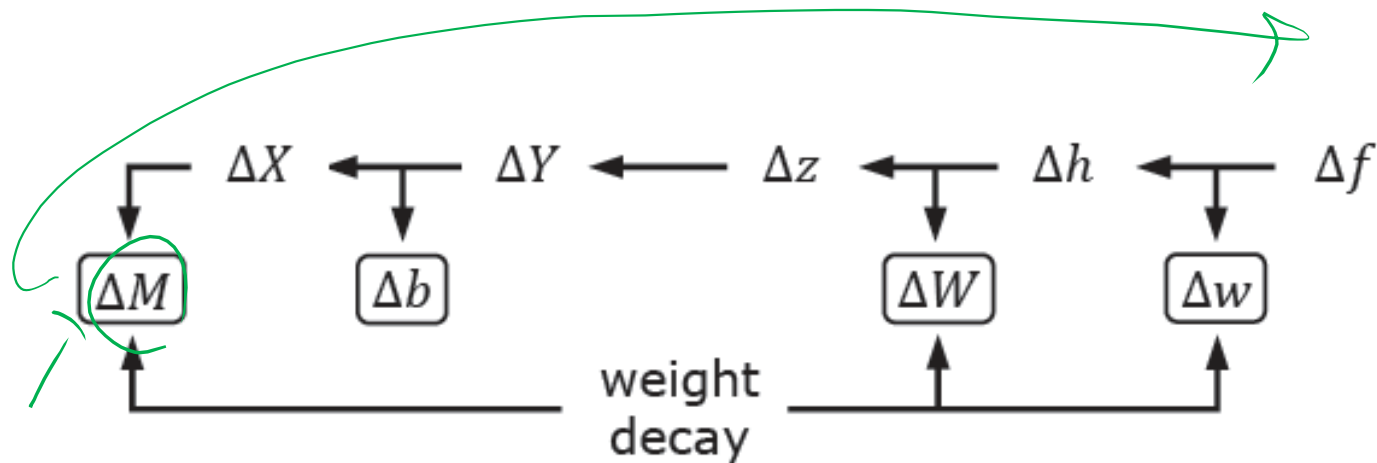
$$\Delta Y_{i,k} = \frac{\Delta z_{2k+1}}{n+m-1} + \begin{cases} \Delta z_{2k+0} & \text{if } i = \operatorname{argmax}(Y_{1,k}, \dots, Y_{n,k}) \\ 0 & \text{otherwise.} \end{cases}$$



Back-propagation : Rectification



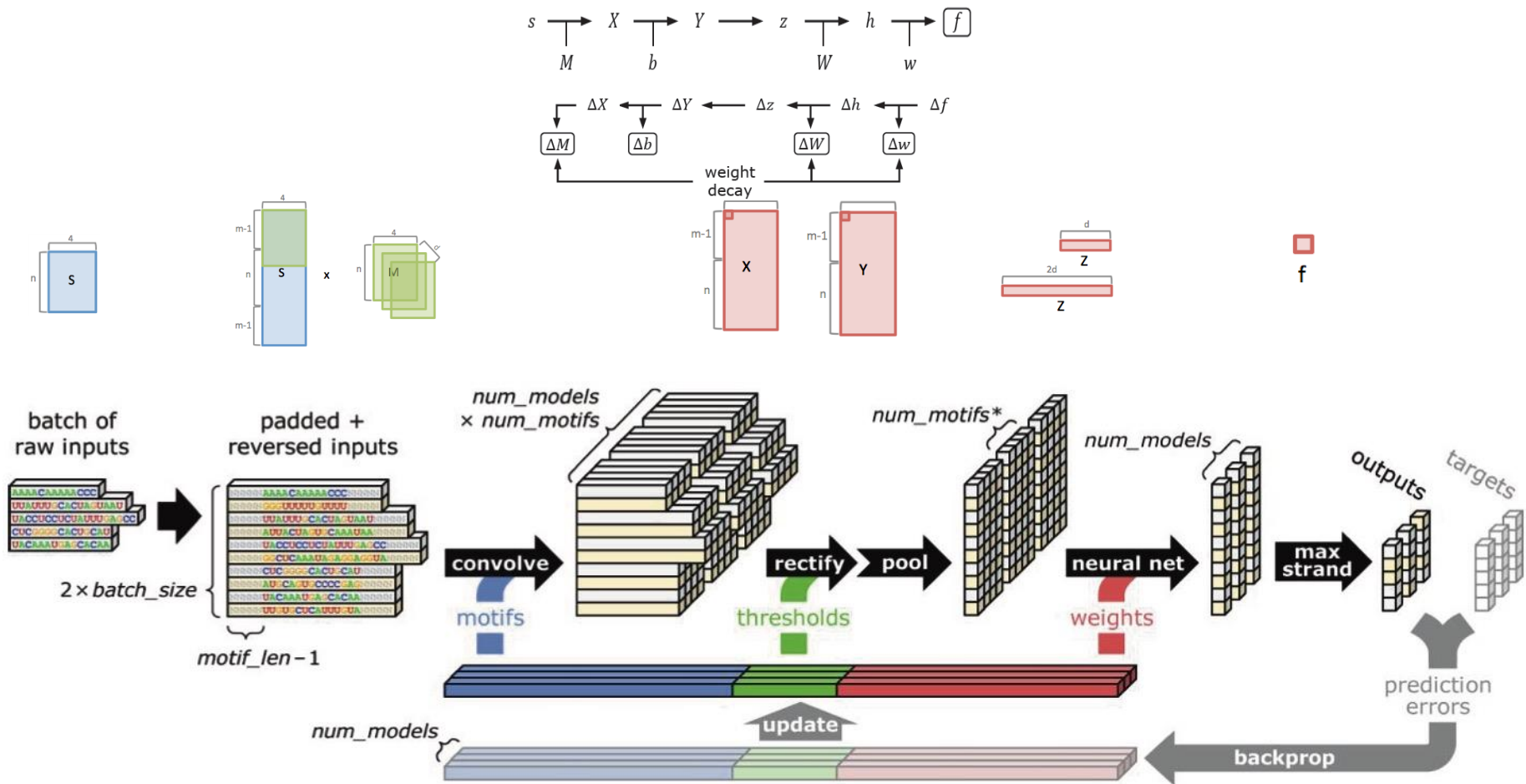
Back-propagation : Convolution



$$\Delta M_{k,j,l} = \sum_{i=1}^{n+m-1} S_{i+j,l} \Delta X_{i,k}$$

Updates

Forward/Backward Propagation



batch_size=5, motif_len=6, num_motifs=4, num_models=3

Summary of DeepBind

- Advantages

- Can be applied to both PBM and ChIP-seq
- Can learn from millions of sequences
- Generalizes well across technologies, without need for correcting for technology-specific biases
- Tolerates noise
- Trains predictive models fully automatically

- Disadvantages

- Model can overfit data (regularizers such as dropout, weight decay, and early stopping are used)
- Parameters do not have any physical meaning

Q&A