# Training example

Language modelling

# Batch-ification

abcdefghijklmnopqrstuvwxyz
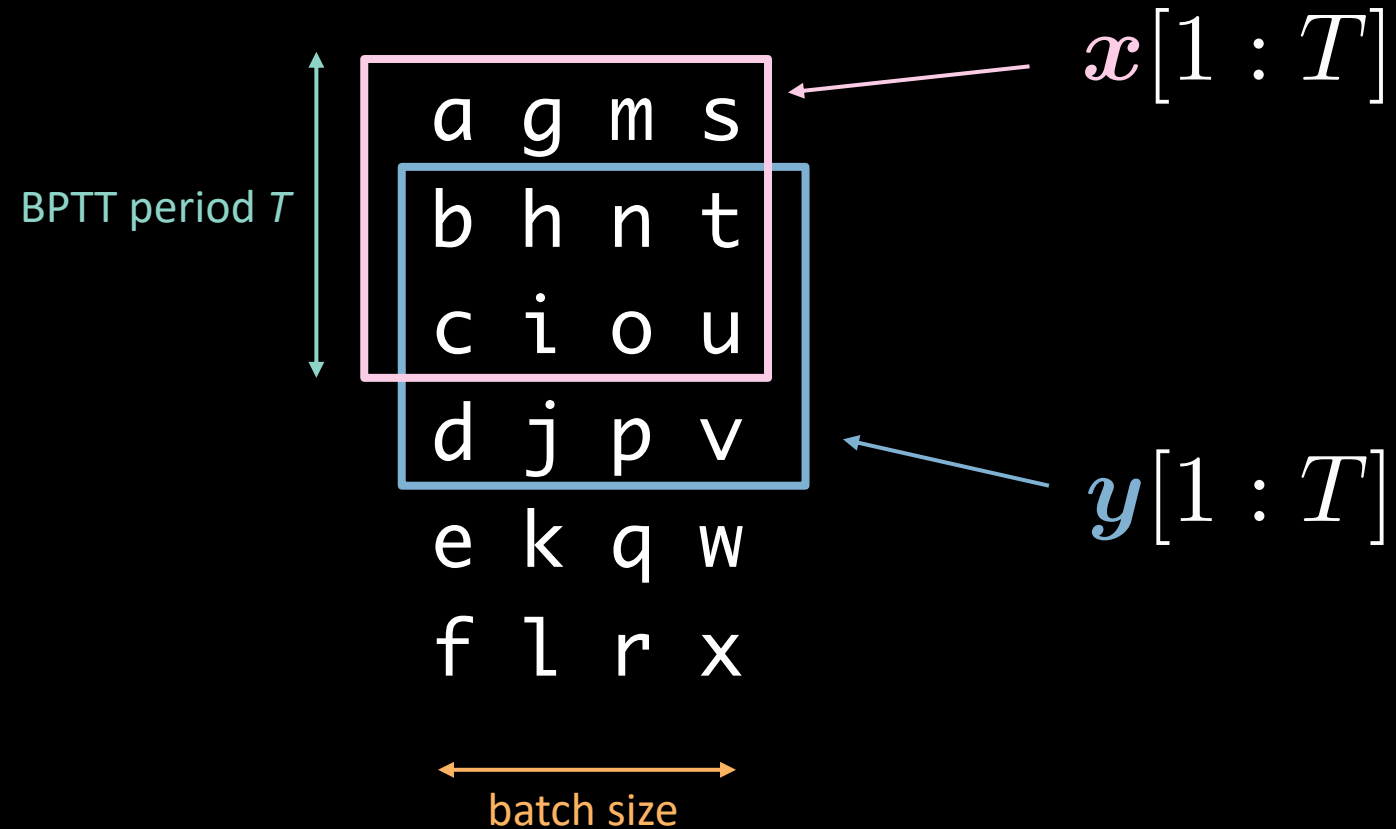
⇩

```
a g m s
b h n t
c i o u
d j p v
e k q w
f l r x
```

↔ batch size

# Get batch (I)

$$x[1:T]$$

$$y[1:T]$$

BPTT period *T*

a g m s
b h n t
c i o u
d j p v
e k q w
f l r x

batch size

# Get batch (II)



Check word_language_model @ github.com/pytorch/examples/

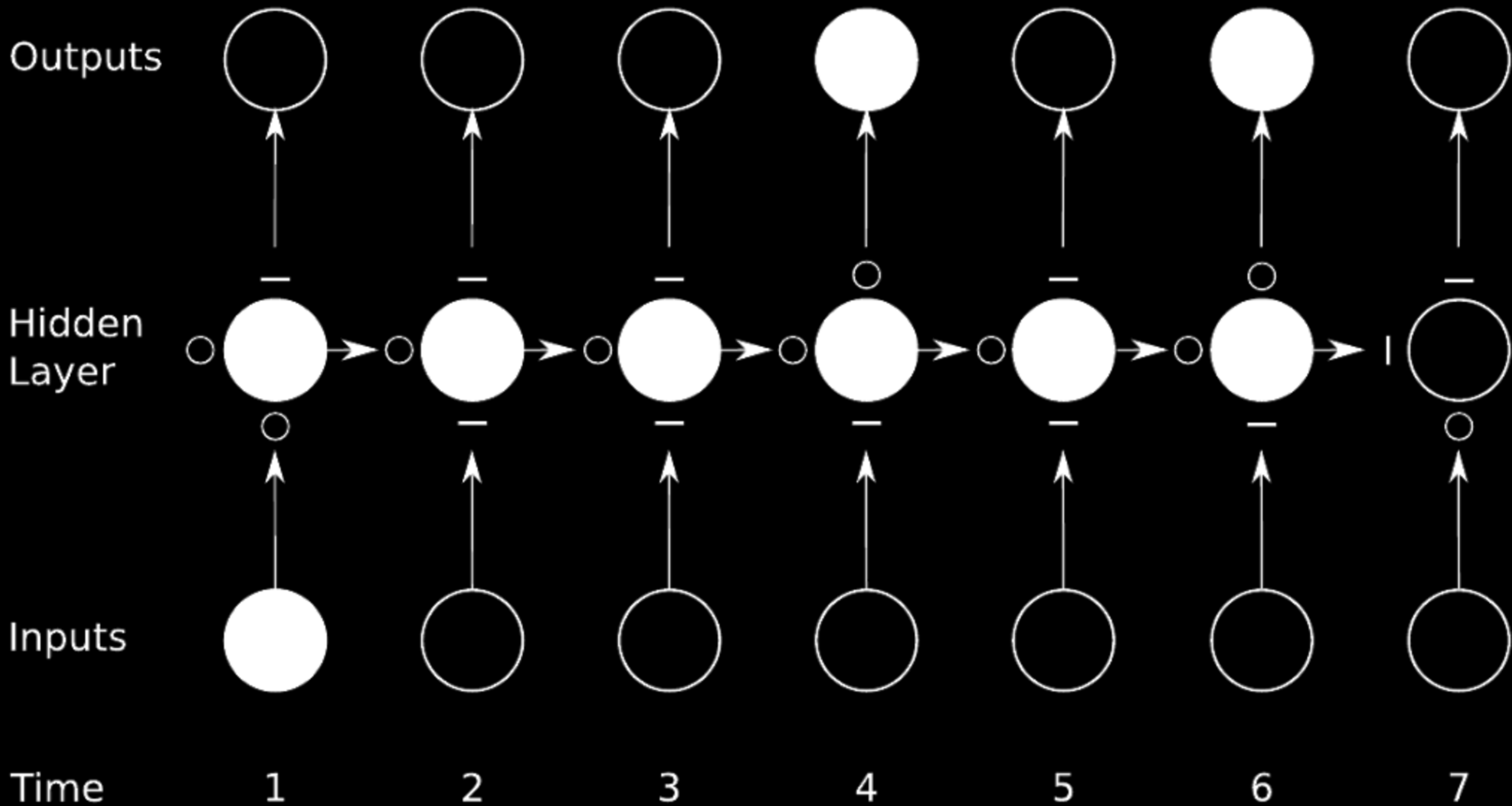# Vanishing & exploding gradients

Limitations of temporally deep nets

Graves (2012) Supervised sequence labelling

Graves (2012) Supervised sequence labelling

# Long Short-Term Memory

Gated RNN

# Controlling the output - OFF

Saturated sigmoid

● = 1
● = 0

$$i[t] = \sigma\left(W_i\begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_i\right)$$

$$f[t] = \sigma\left(W_f\begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_f\right)$$

$$o[t] = \sigma\left(W_o\begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_o\right)$$

$$\tilde{c}[t] = \tanh\left(W_c\begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_c\right)$$

$$c[t] = f[t] \odot c[t-1] + i[t] \odot \tilde{c}[t]$$

$$h[t] = o[t] \odot \tanh(c[t])$$

# Controlling the output - ON

Saturated sigmoid
- = 1
- = 0

$$i[t] = \sigma\left(W_i \begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_i\right)$$

$$f[t] = \sigma\left(W_f \begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_f\right)$$

$$o[t] = \sigma\left(W_o \begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_o\right)$$

$$\tilde{c}[t] = \tanh\left(W_c \begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_c\right)$$

$$c[t] = f[t] \odot c[t-1] + i[t] \odot \tilde{c}[t]$$

$$h[t] = o[t] \odot \tanh(c[t])$$

# Controlling the memory - reset

Saturated sigmoid
- = 1
- = 0

$$i[t] = \sigma\left(W_i \begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_i\right)$$

$$f[t] = \sigma\left(W_f \begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_f\right)$$

$$o[t] = \sigma\left(W_o \begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_o\right)$$

$$\tilde{c}[t] = \tanh\left(W_c \begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_c\right)$$

$$c[t] = f[t] \odot c[t-1] + i[t] \odot \tilde{c}[t]$$

$$h[t] = o[t] \odot \tanh(c[t])$$

# Controlling the memory - keep

Saturated sigmoid
- 🟢 = 1
- 🔴 = 0

$$i[t] = \sigma\left(W_i\left[\begin{smallmatrix}x[t]\\h[t-1]\end{smallmatrix}\right] + b_i\right)$$

$$f[t] = \sigma\left(W_f\left[\begin{smallmatrix}x[t]\\h[t-1]\end{smallmatrix}\right] + b_f\right)$$

$$o[t] = \sigma\left(W_o\left[\begin{smallmatrix}x[t]\\h[t-1]\end{smallmatrix}\right] + b_o\right)$$

$$\tilde{c}[t] = \tanh\left(W_c\left[\begin{smallmatrix}x[t]\\h[t-1]\end{smallmatrix}\right] + b_c\right)$$

$$c[t] = f[t] \odot c[t-1] + i[t] \odot \tilde{c}[t]$$

$$h[t] = o[t] \odot \tanh(c[t])$$

# Controlling the memory - write

Saturated sigmoid
- 🟢 = 1
- 🔴 = 0

$$i[t] = \sigma\big(\boldsymbol{W}_i\big[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\big] + \boldsymbol{b}_i\big)$$

$$\boldsymbol{f}[t] = \sigma\big(\boldsymbol{W}_f\big[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\big] + \boldsymbol{b}_f\big)$$

$$\boldsymbol{o}[t] = \sigma\big(\boldsymbol{W}_o\big[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\big] + \boldsymbol{b}_o\big)$$

$$\tilde{\boldsymbol{c}}[t] = \tanh\big(\boldsymbol{W}_c\big[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\big] + \boldsymbol{b}_c\big)$$

$$\boldsymbol{c}[t] = \boldsymbol{f}[t] \odot \boldsymbol{c}[t-1] + \boldsymbol{i}[t] \odot \tilde{\boldsymbol{c}}[t]$$

$$\boldsymbol{h}[t] = \boldsymbol{o}[t] \odot \tanh(\boldsymbol{c}[t])$$