# Wyprawa w głąb LSTMa

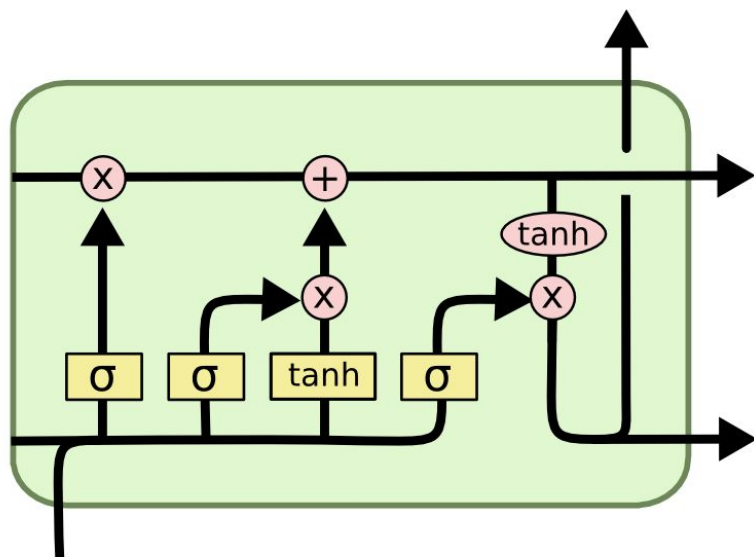CORE, 13 maja 2020
Julia Bazińska

# Long Short-Term Memory

*"Recurrent networks can in principle use their feedback connections to store representations of recent input events in the form of activations ("short-term memory", as opposed to "long-term memory embodied by slowly changing weights)."*

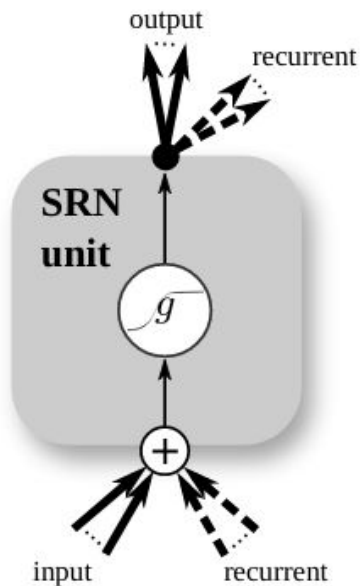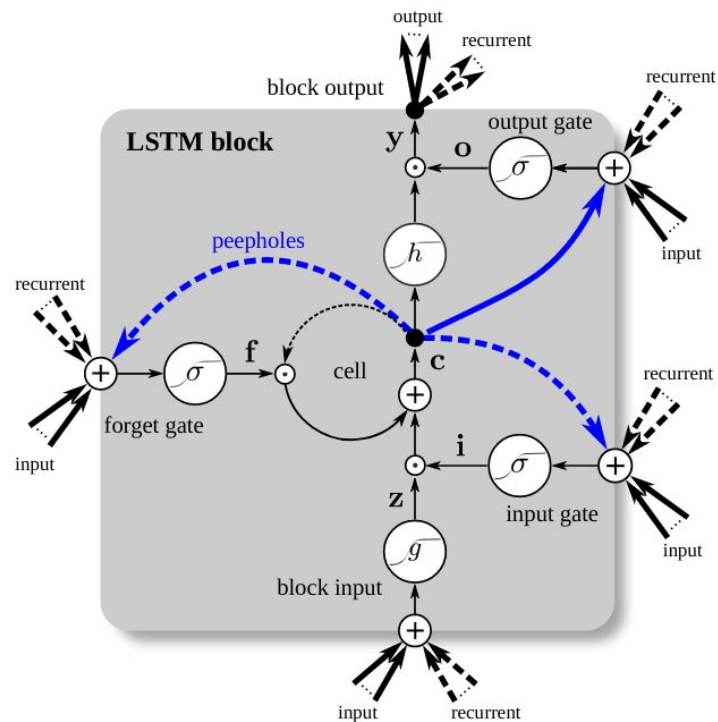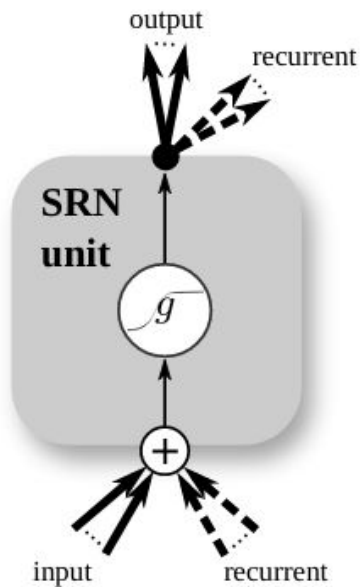*Hochreiter et al, "Long Short-Term Memory", 1997*

# Agenda

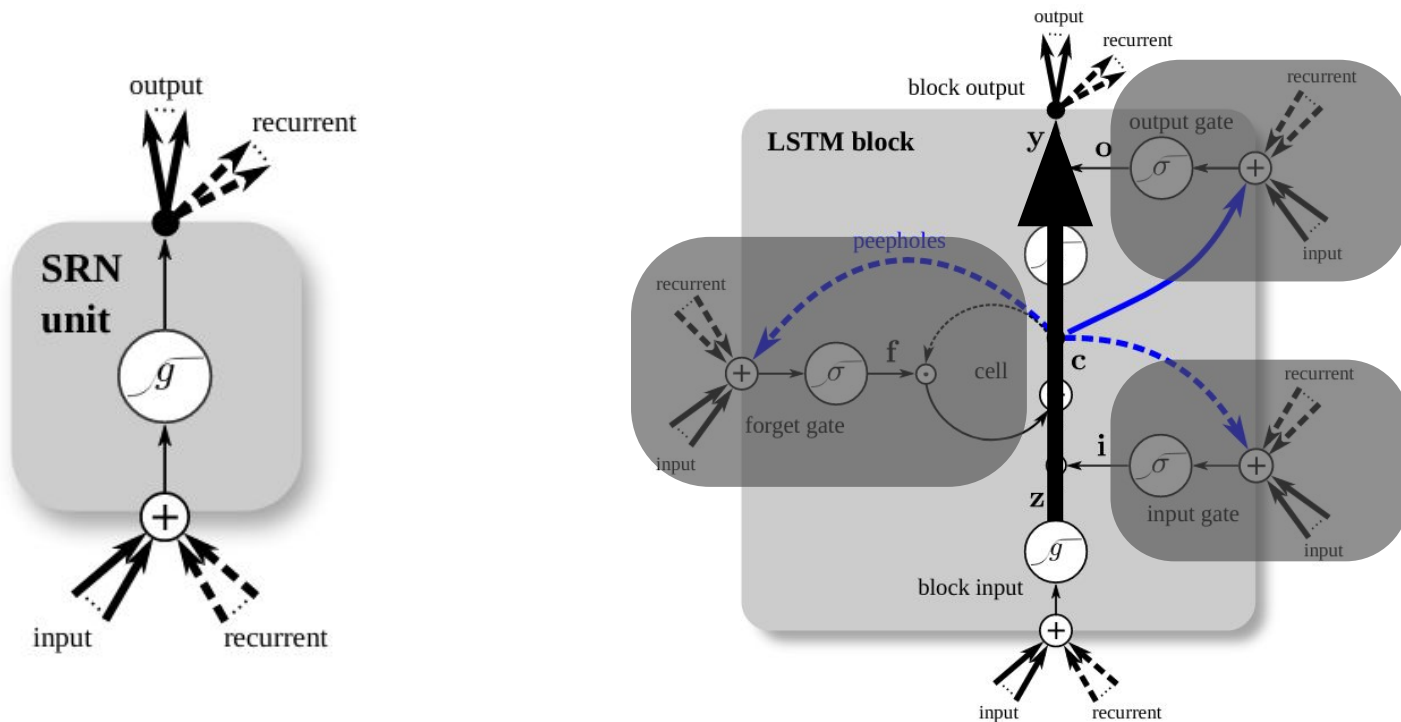# LSTM vs prosta jednostka rekurencyjna

# Najprostsza jednostka rekurencyjna



$$y_t = g(W\,x_t + R\,y_{t-1} + b)$$

# Czy to na pewno jest potrzebne?

# Co tam się właściwie dzieje?

# Co tam się właściwie dzieje?



$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$
$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad \qquad \textit{block input}$$
$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$
$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad \qquad \textit{input gate}$$
$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$
$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad \qquad \textit{forget gate}$$
$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad \qquad \textit{cell}$$
$$\bar{\mathbf{o}}^t = \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \mathbf{p}_o \odot \mathbf{c}^t + \mathbf{b}_o$$
$$\mathbf{o}^t = \sigma(\bar{\mathbf{o}}^t) \qquad \qquad \textit{output gate}$$
$$\mathbf{y}^t = h(\mathbf{c}^t) \odot \mathbf{o}^t \qquad \qquad \textit{block output}$$

# Co tam się właściwie dzieje?



$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$

$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad\qquad\qquad \textit{block input}$$

$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$

$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad\qquad\qquad \textit{input gate}$$

$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$

$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad\qquad\qquad \textit{forget gate}$$

$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad\qquad \textit{cell}$$

$$\bar{\mathbf{o}}^t = \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \mathbf{p}_o \odot \mathbf{c}^t + \mathbf{b}_o$$

$$\mathbf{o}^t = \sigma(\bar{\mathbf{o}}^t) \qquad\qquad\qquad \textit{output gate}$$

$$\mathbf{y}^t = h(\mathbf{c}^t) \odot \mathbf{o}^t \qquad\qquad\qquad \textit{block output}$$

# Co tam się właściwie dzieje?



$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$
$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t)$$

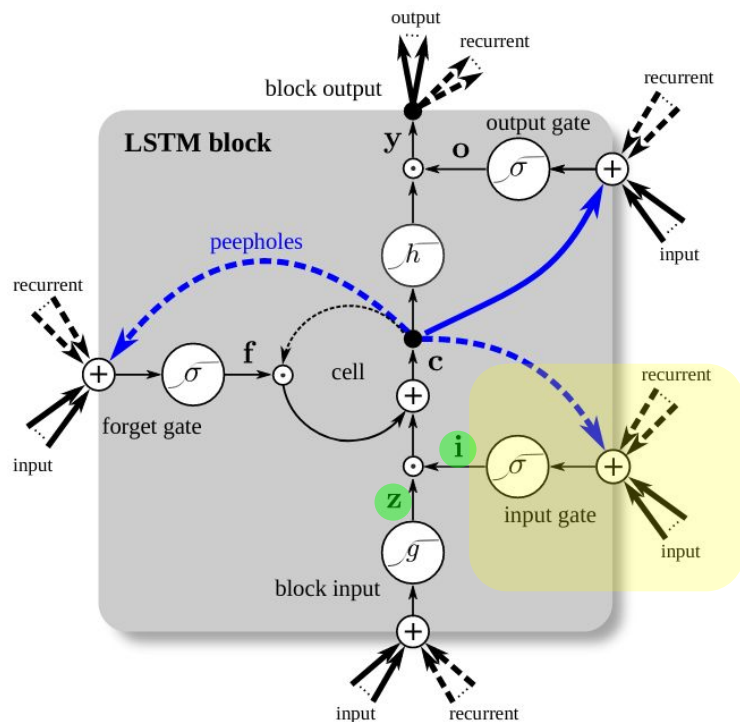*block input*

# Co tam się właściwie dzieje?



$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$
$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad\qquad\qquad\qquad \textit{block input}$$
$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$
$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad\qquad\qquad\qquad \textit{input gate}$$

# Co tam się właściwie dzieje?



$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$

$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad\qquad\qquad block\ input$$
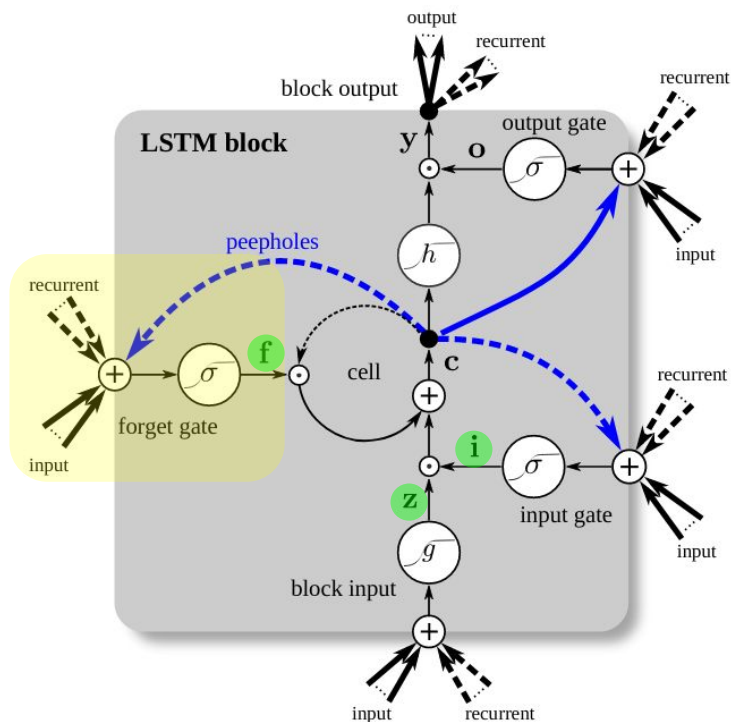
$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$

$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad\qquad\qquad input\ gate$$

$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$

$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad\qquad\qquad forget\ gate$$

# Co tam się właściwie dzieje?



$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$
$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad\qquad block\ input$$
$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$
$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad\qquad input\ gate$$
$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$
$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad\qquad forget\ gate$$
$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad\qquad cell$$

# Co tam się właściwie dzieje?



$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$

$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad \textit{block input}$$

$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$

$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad \textit{input gate}$$

$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$

$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad \textit{forget gate}$$

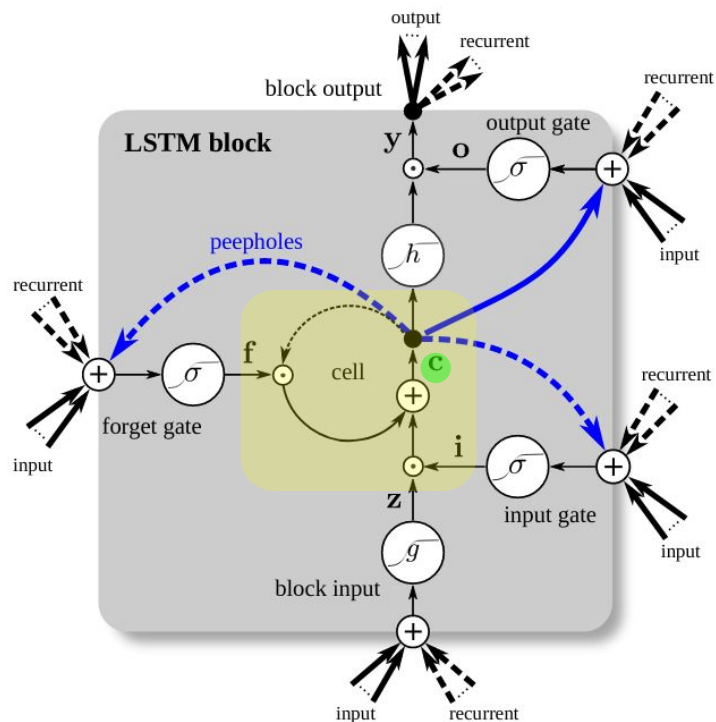$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad \textit{cell}$$

# Co tam się właściwie dzieje?



$$\bar{\mathbf{z}}^t = \mathbf{W}_z\mathbf{x}^t + \mathbf{R}_z\mathbf{y}^{t-1} + \mathbf{b}_z$$

$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad \textit{block input}$$

$$\bar{\mathbf{i}}^t = \mathbf{W}_i\mathbf{x}^t + \mathbf{R}_i\mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$

$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad \textit{input gate}$$

$$\bar{\mathbf{f}}^t = \mathbf{W}_f\mathbf{x}^t + \mathbf{R}_f\mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$

$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad \textit{forget gate}$$

$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad \textit{cell}$$

# Co tam się właściwie dzieje?



$$\bar{\mathbf{z}}^t = \mathbf{W}_z\mathbf{x}^t + \mathbf{R}_z\mathbf{y}^{t-1} + \mathbf{b}_z$$

$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad\qquad\qquad\qquad block\ input$$

$$\bar{\mathbf{i}}^t = \mathbf{W}_i\mathbf{x}^t + \mathbf{R}_i\mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$

$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad\qquad\qquad\qquad input\ gate$$

$$\bar{\mathbf{f}}^t = \mathbf{W}_f\mathbf{x}^t + \mathbf{R}_f\mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$
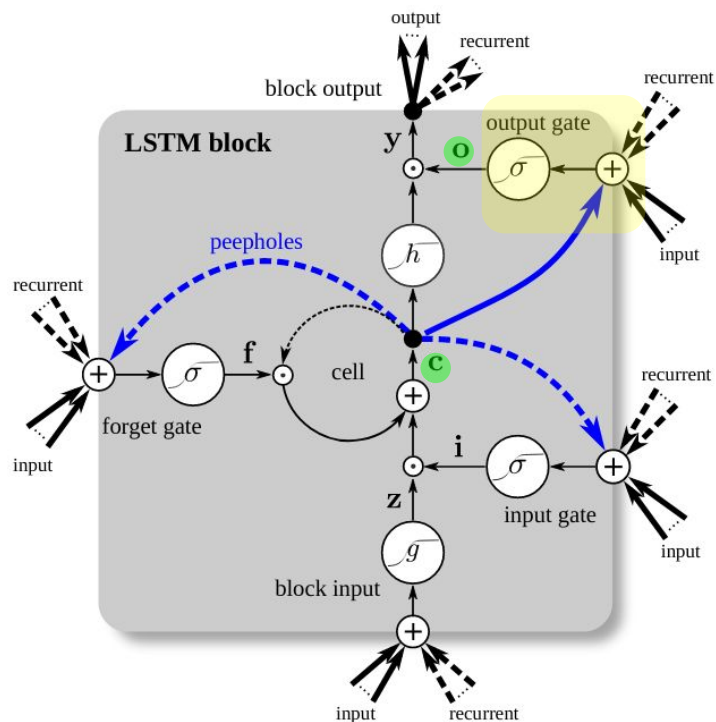
$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad\qquad\qquad\qquad forget\ gate$$

$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad\qquad cell$$

# Co tam się właściwie dzieje?



$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$
$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad \qquad \textit{block input}$$
$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$
$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad \qquad \textit{input gate}$$
$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$
$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad \qquad \textit{forget gate}$$
$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad \qquad \textit{cell}$$
$$\bar{\mathbf{o}}^t = \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \mathbf{p}_o \odot \mathbf{c}^t + \mathbf{b}_o$$
$$\mathbf{o}^t = \sigma(\bar{\mathbf{o}}^t) \qquad \qquad \textit{output gate}$$

# Co tam się właściwie dzieje?



$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$

$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad \textit{block input}$$

$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$

$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad \textit{input gate}$$

$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$

$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad \textit{forget gate}$$

$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad \textit{cell}$$

$$\bar{\mathbf{o}}^t = \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \mathbf{p}_o \odot \mathbf{c}^t + \mathbf{b}_o$$

$$\mathbf{o}^t = \sigma(\bar{\mathbf{o}}^t) \qquad \textit{output gate}$$

$$\mathbf{y}^t = h(\mathbf{c}^t) \odot \mathbf{o}^t \qquad \textit{block output}$$

# Co tam się właściwie dzieje?



$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$

$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad \qquad \textit{block input}$$

$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$

$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad \qquad \textit{input gate}$$

$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$

$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad \qquad \textit{forget gate}$$

$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad \qquad \textit{cell}$$

$$\bar{\mathbf{o}}^t = \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \mathbf{p}_o \odot \mathbf{c}^t + \mathbf{b}_o$$

$$\mathbf{o}^t = \sigma(\bar{\mathbf{o}}^t) \qquad \qquad \textit{output gate}$$

$$\mathbf{y}^t = h(\mathbf{c}^t) \odot \mathbf{o}^t \qquad \qquad \textit{block output}$$

# Peephole connections



$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$

$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad \textit{block input}$$

$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \boxed{\mathbf{p}_i \odot \mathbf{c}^{t-1}} + \mathbf{b}_i$$

$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad \textit{input gate}$$

$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \boxed{\mathbf{p}_f \odot \mathbf{c}^{t-1}} + \mathbf{b}_f$$

$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad \textit{forget gate}$$

$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad \textit{cell}$$

$$\bar{\mathbf{o}}^t = \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \boxed{\mathbf{p}_o \odot \mathbf{c}^t} + \mathbf{b}_o$$

$$\mathbf{o}^t = \sigma(\bar{\mathbf{o}}^t) \qquad \textit{output gate}$$

$$\mathbf{y}^t = h(\mathbf{c}^t) \odot \mathbf{o}^t \qquad \textit{block output}$$

# LSTM vs prosta jednostka reukrencyjna

# Warianty LSTM

# Warianty

1. Usunięcie gate'a:
   a. **No Input Gate**
   b. No Forget Gate
   c. No Output Gate
2. Usunięcie aktywacji:
   a. No Input Activation Function
   b. No Output Activation Function
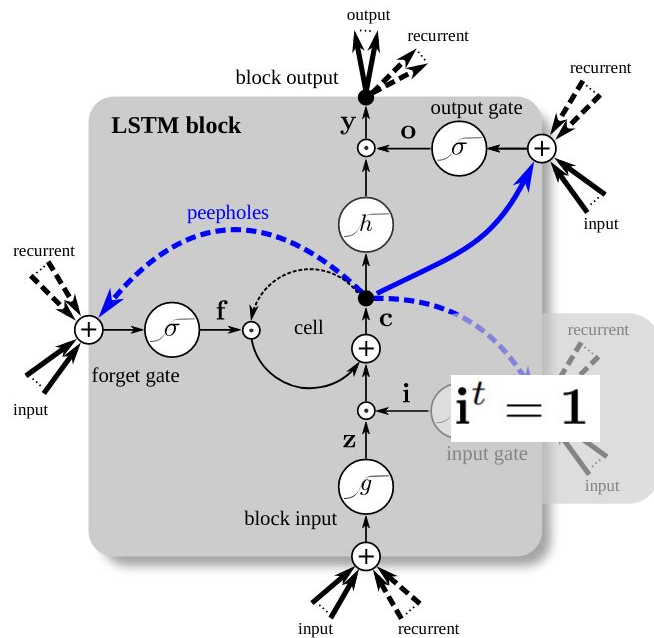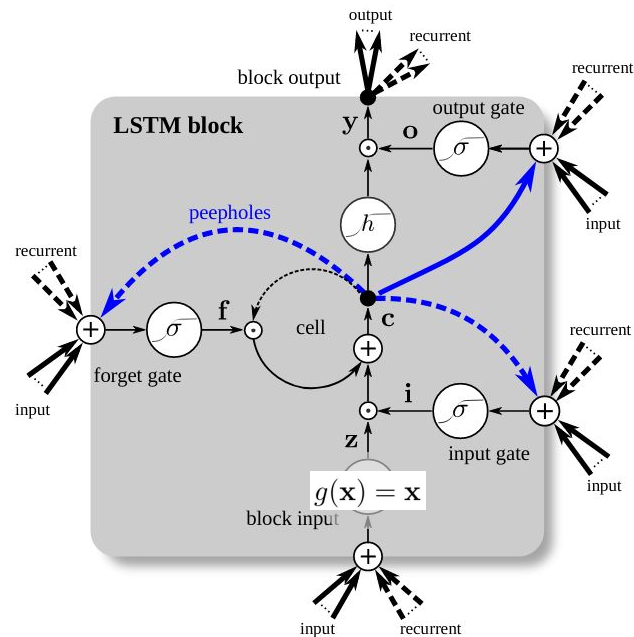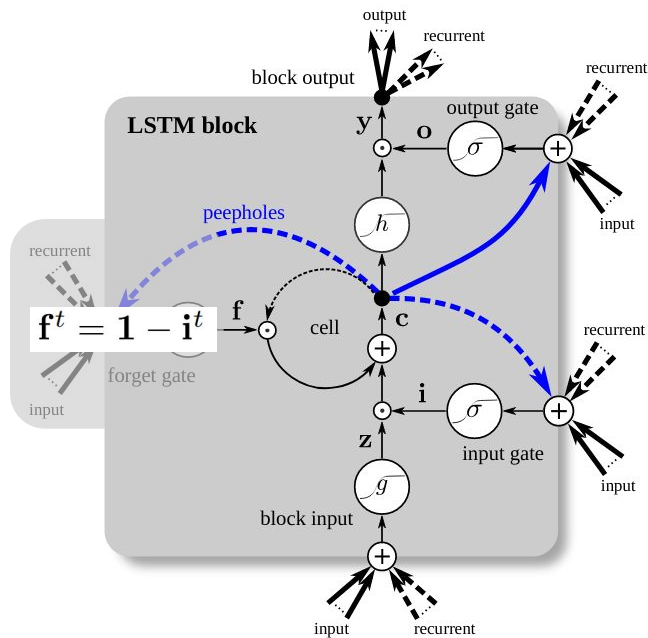3. Couple Input and Forget Gate
4. No Peepholes
5. Full Gate Recurrence

# Warianty

1. Usunięcie gate'a:
   a. **No Input Gate**
   b. No Forget Gate
   c. No Output Gate
2. Usunięcie aktywacji:
   a. No Input Activation Function
   b. No Output Activation Function
3. Couple Input and Forget Gate
4. No Peepholes
5. Full Gate Recurrence
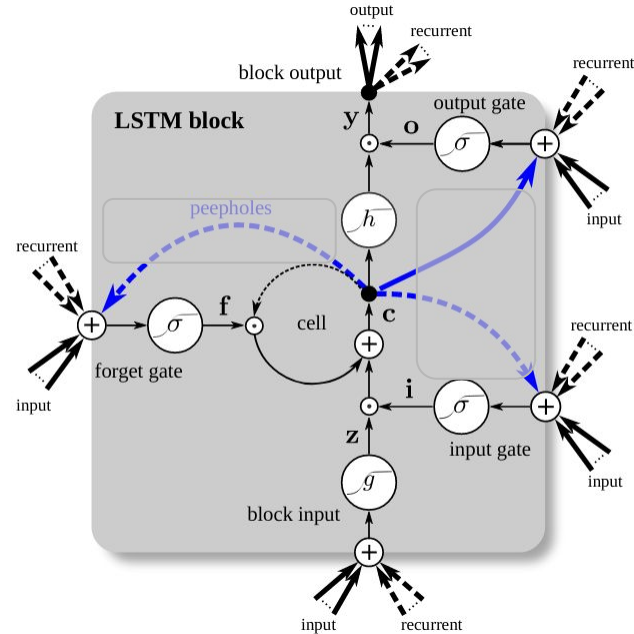
# Warianty

1. Usunięcie gate'a:
   a. **No Input Gate**
   b. No Forget Gate
   c. No Output Gate
2. Usunięcie aktywacji:
   a. No Input Activation Function
   b. No Output Activation Function
3. Couple Input and Forget Gate
4. No Peepholes
5. Full Gate Recurrence

# Warianty

1. Usunięcie gate'a:
   a. No Input Gate
   b. No Forget Gate
   c. No Output Gate
2. Usunięcie aktywacji:
   a. **No Input Activation Function**
   b. No Output Activation Function
3. Couple Input and Forget Gate
4. No Peepholes
5. Full Gate Recurrence
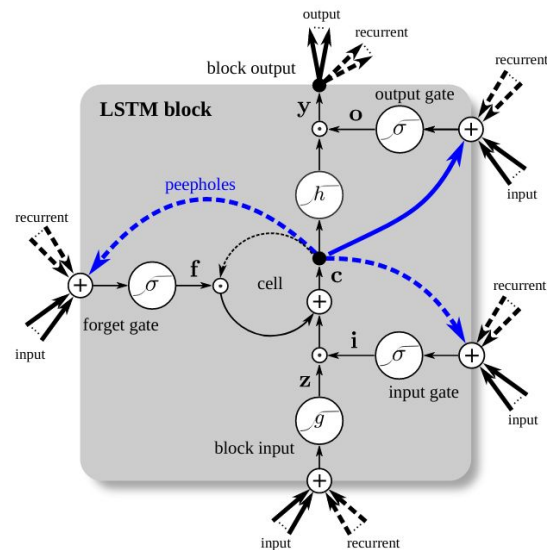
# Warianty

1. Usunięcie gate'a:
   a. No Input Gate
   b. No Forget Gate
   c. No Output Gate
2. Usunięcie aktywacji:
   a. No Input Activation Function
   b. No Output Activation Function
3. **Coupled Input and Forget Gate**
4. No Peepholes
5. Full Gate Recurrence

# Warianty

1. Usunięcie gate'a:
   a. No Input Gate
   b. No Forget Gate
   c. No Output Gate
2. Usunięcie aktywacji:
   a. No Input Activation Function
   b. No Output Activation Function
3. Coupled Input and Forget Gate
4. **No Peepholes**
5. Full Gate Recurrence

# Warianty

1. Usunięcie gate'a:
   a. No Input Gate
   b. No Forget Gate
   c. No Output Gate
2. Usunięcie aktywacji:
   a. No Input Activation Function
   b. No Output Activation Function
3. Coupled Input and Forget Gate
4. No Peepholes
5. **Full Gate Recurrence**

   **każdy gate + wyniki z poprzedniego kroku ze wszystkich gate'ów**

$$+ \mathbf{R}_{ii}\mathbf{i}^{t-1} + \mathbf{R}_{fi}\mathbf{f}^{t-1} + \mathbf{R}_{oi}\mathbf{o}^{t-1}$$



$$\mathbf{i}^t = \mathbf{W}_i\mathbf{x}^t + \mathbf{R}_i\mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$
$$+ \mathbf{R}_{ii}\mathbf{i}^{t-1} + \mathbf{R}_{fi}\mathbf{f}^{t-1} + \mathbf{R}_{oi}\mathbf{o}^{t-1}$$
$$\bar{\mathbf{f}}^t = \mathbf{W}_f\mathbf{x}^t + \mathbf{R}_f\mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$
$$+ \mathbf{R}_{if}\mathbf{i}^{t-1} + \mathbf{R}_{ff}\mathbf{f}^{t-1} + \mathbf{R}_{of}\mathbf{o}^{t-1}$$
$$\bar{\mathbf{o}}^t = \mathbf{W}_o\mathbf{x}^t + \mathbf{R}_o\mathbf{y}^{t-1} + \mathbf{p}_o \odot \mathbf{c}^{t-1} + \mathbf{b}_o$$
$$+ \mathbf{R}_{io}\mathbf{i}^{t-1} + \mathbf{R}_{fo}\mathbf{f}^{t-1} + \mathbf{R}_{oo}\mathbf{o}^{t-1}$$

# Trzy benchmarki

# Benchmarki / architektury

- TIMIT - rozpoznawanie mowy (fonetyczne)
  - 1 warstwa, dwukierunkowy LSTM
- IAM Online - rozpoznawanie pisma
  - 1 warstwa, dwukierunkowy LSTM
- JSB Chorales - predykcja muzyki
  - 1 warstwa, jednokierunkowy LSTM

# Eksperymenty i wyniki

# Eksperymenty

- każdy wariant x każdy benchmark
- random search hiperparametrów:
  - rozmiar warstwy
  - szum na wejściu
  - momentum
  - learning rate



Kula Zorba tocząca się w dół zbocza, obrazująca dążenie do minimum funkcji straty, *iStock*.
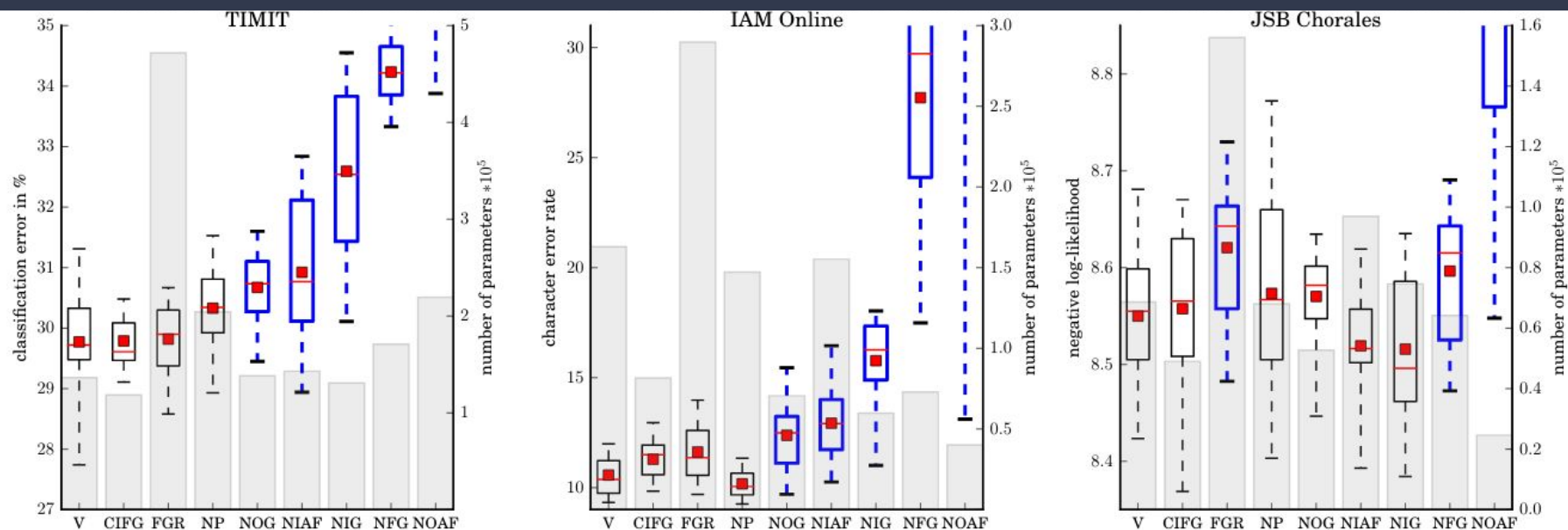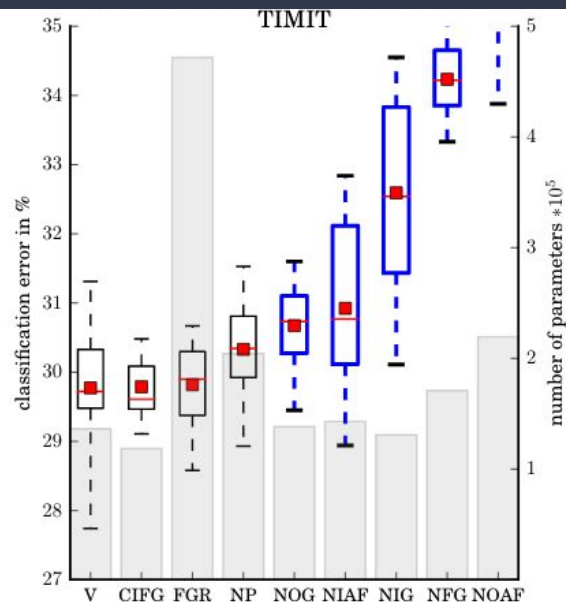
# Wyniki



Figure 3. *Test set* performance for all 200 trials (top) and for the best 10% (bottom) trials (according to the *validation set*) for each dataset and variant. Boxes show the range between the 25<sup>th</sup> and the 75<sup>th</sup> percentile of the data, while the whiskers indicate the whole range. The red dot represents the mean and the red line the median of the data. The boxes of variants that differ significantly from the vanilla LSTM are shown in blue with thick lines. The grey histogram in the background presents the average number of parameters for the top 10% performers of every variant.
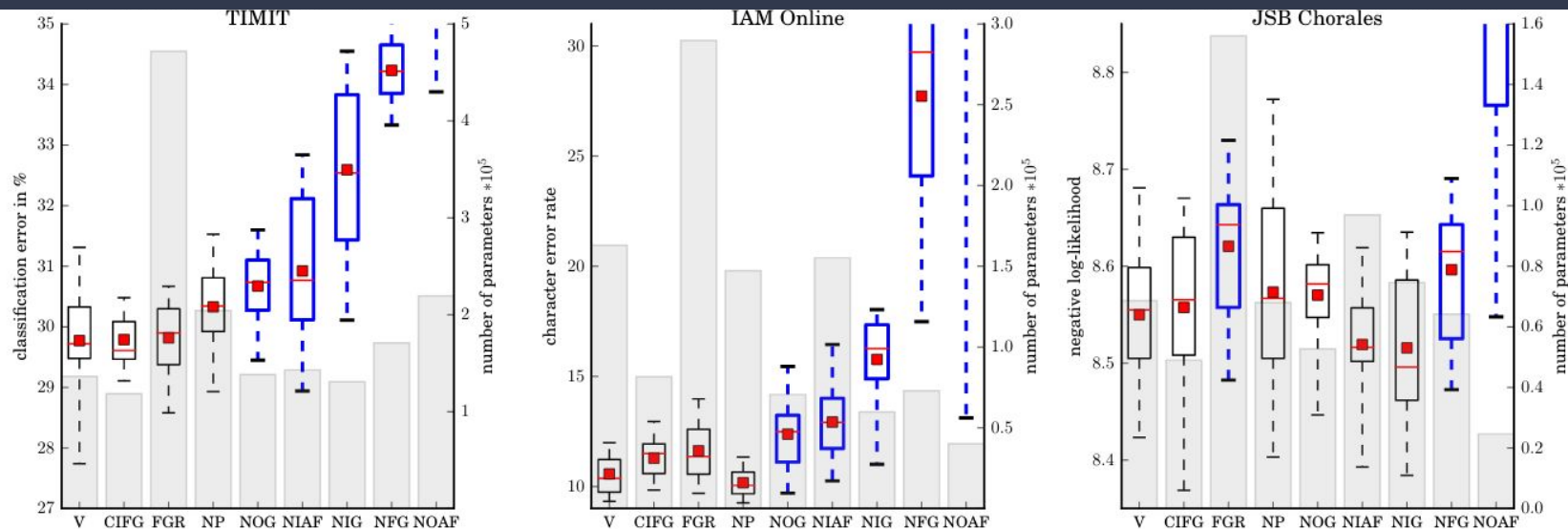
# Wyniki



V – Vanilla,
CIFG – Coupled Input Forget Gates,
FGR – Full Gate Recurrence,
NP – No Peepholes,
NOG – No Output Gate,
NIAF – No Input Activation Function,
NIG – No Input Gate,
NFG – No Forget Gate,
NOAF – No Output Activation Function

- Output Activation, Forget Gate – bardzo ważne!
- …ale już CIFG nadal radzi sobie dobrze
- NP – nieznaczne zmiany w wynikach…
- FGR nie pomaga, ale dodaje parametrów
- NIG, NOG, NIAF – gorzej w rozpoznawaniu pisma i mowy, ale dla muzyki ok.

Figure 3. *Test set* performance for all 200 trials (top) and for the best 10% (bottom) trials (according to the *validation set*) for each dataset and variant. Boxes show the range between the 25th and the 75th percentile of the data, while the whiskers indicate the whole range. The red dot represents the mean and the red line the median of the data. The boxes of variants that differ significantly from the vanilla LSTM are shown in blue with thick lines. The grey histogram in the background presents the average number of parameters for the top 10% performers of every variant.

# Wyniki



TIMIT

IAM Online

JSB Chorales

- Output Activation, Forget Gate – bardzo ważne!
- …ale już CIFG nadal radzi sobie dobrze
- NP – nieznaczne zmiany w wynikach…
- FGR nie pomaga, ale dodaje parametrów
- NIG, NOG, NIAF – gorzej w rozpoznawaniu pisma i mowy, ale dla muzyki ok.

V – Vanilla,
CIFG – Coupled Input Forget Gates,
FGR – Full Gate Recurrence,
NP – No Peepholes,
NOG – No Output Gate,
NIAF – No Input Activation Function,
NIG – No Input Gate,
NFG – No Forget Gate,
NOAF – No Output Activation Function