

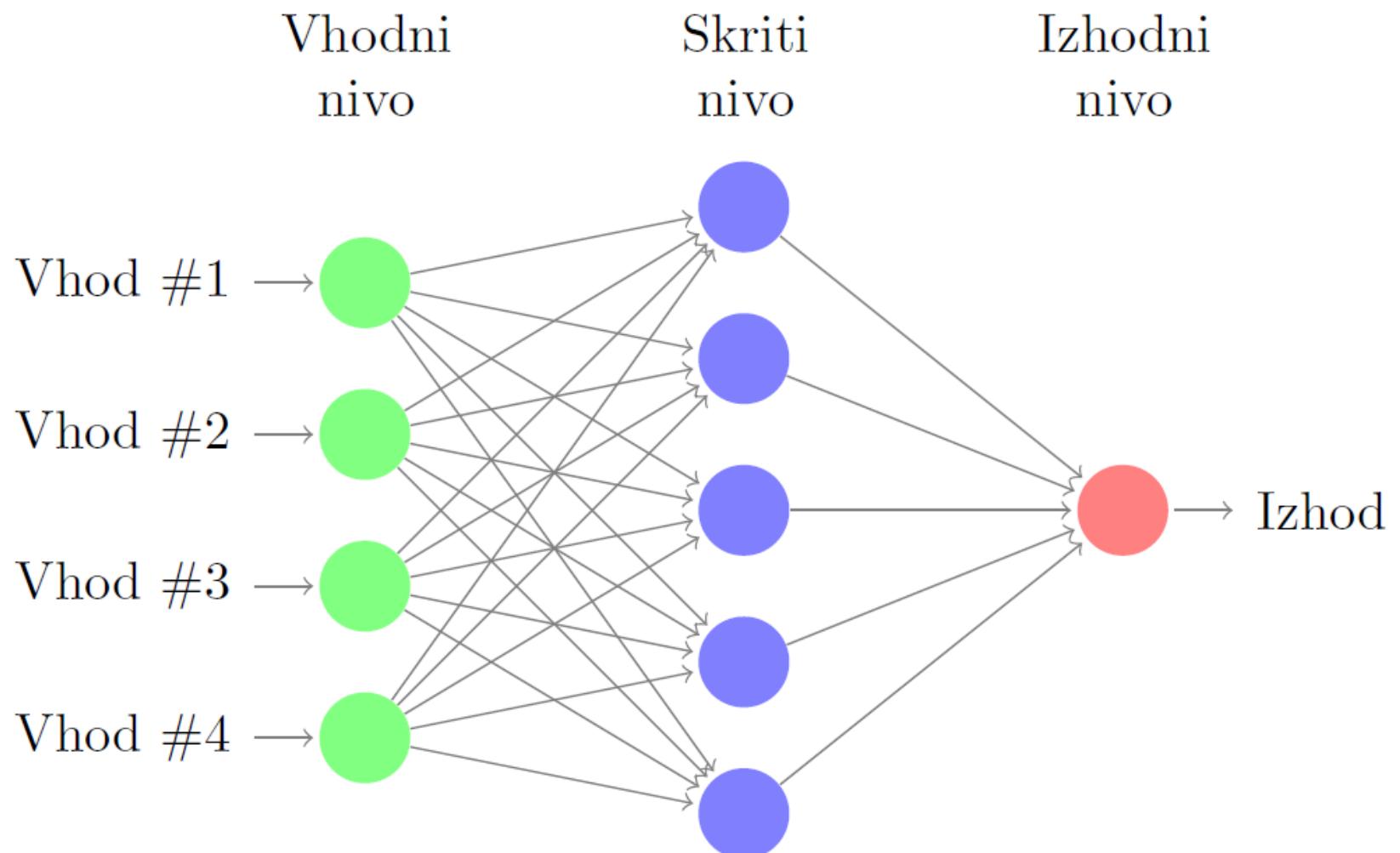
# Uporaba globokih konvolucijskih nevronske mrež na surovem besedilu

Žiga Pušnik

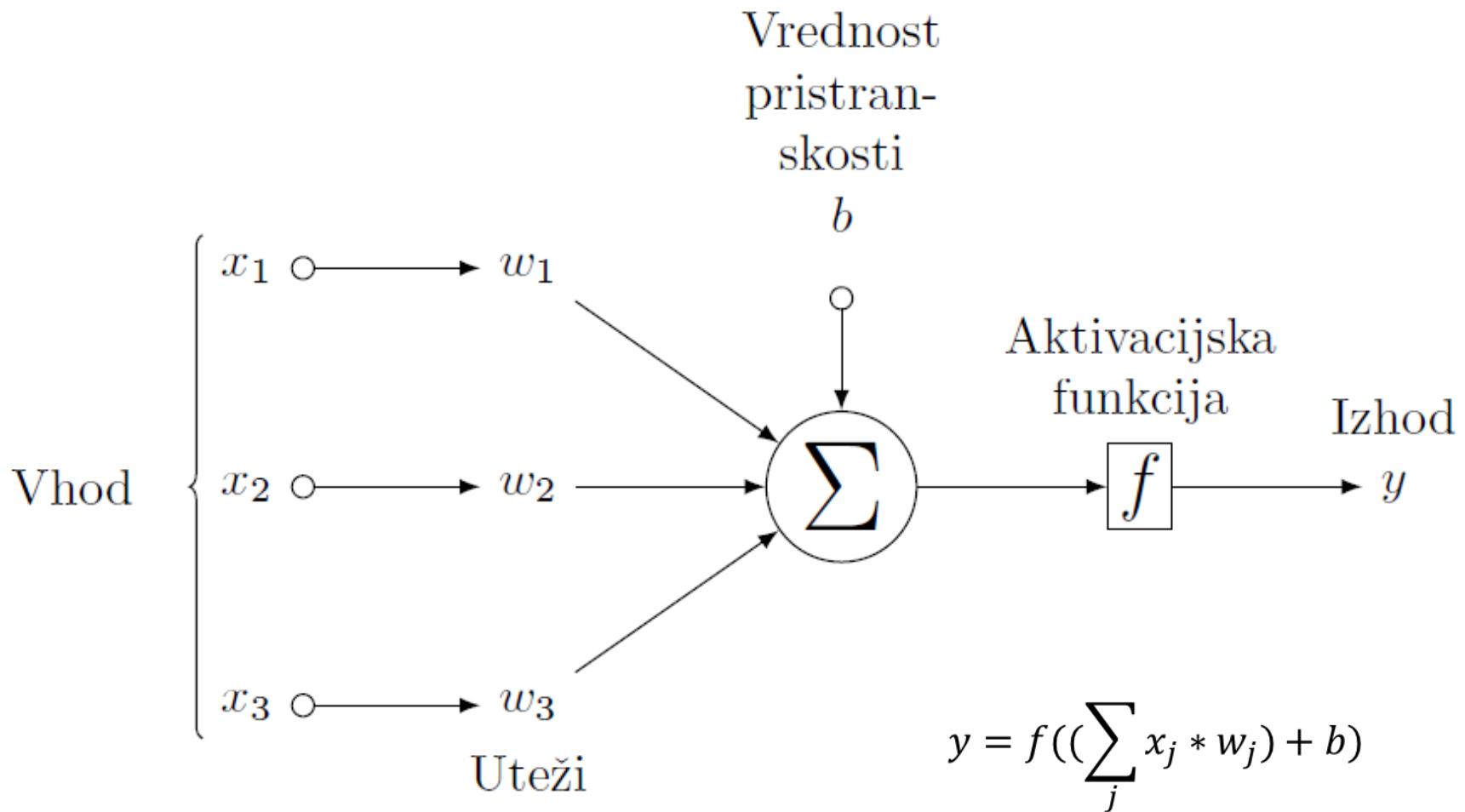
- Ali lahko računalnik razume besedilo s konvolucijskimi nevronskimi mrežami?
- *Text understanding from scratch*
  - (98% – 99%) klasifikacijska točnost
- klasifikacija člankov
- predprocesiranje besedila pri obdelavi naravnega jezika

- vzporednica z možgani
- aproksimacija poljubne funkcije (teoretično)
- globje arhitekture za kompleksnejše funkcije
- težave pri globokih arhitekturah

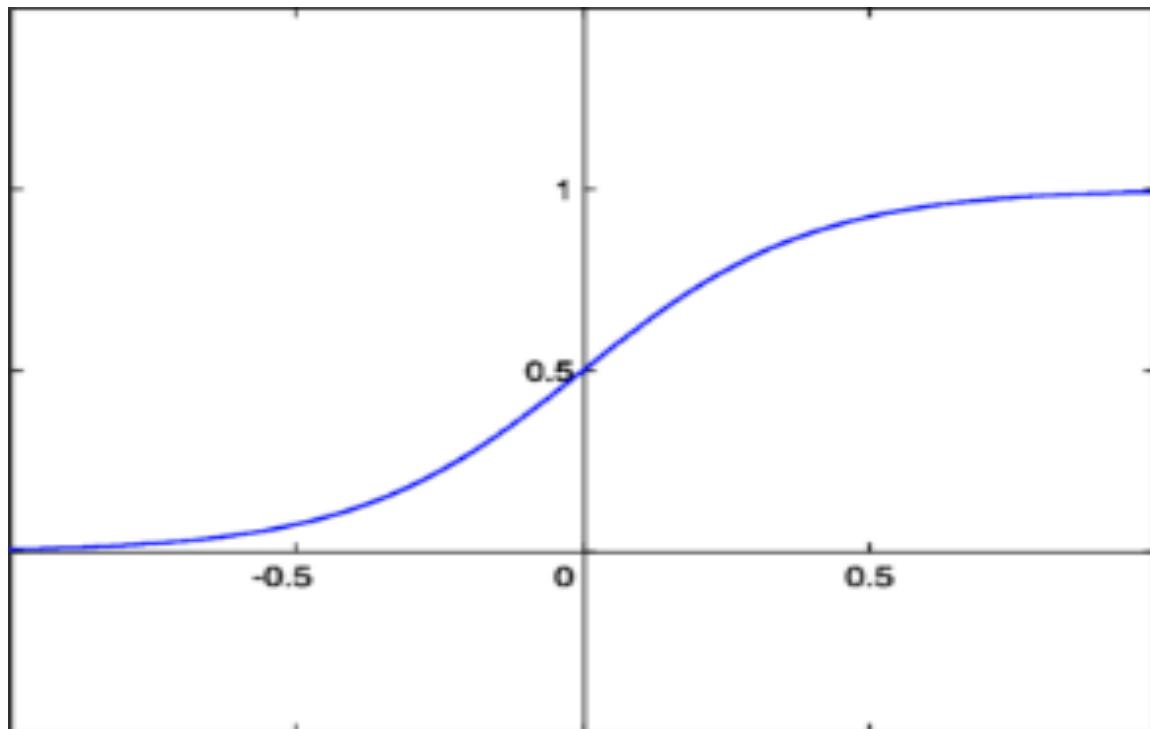
# Nevronske mreže



# Shema nevrona



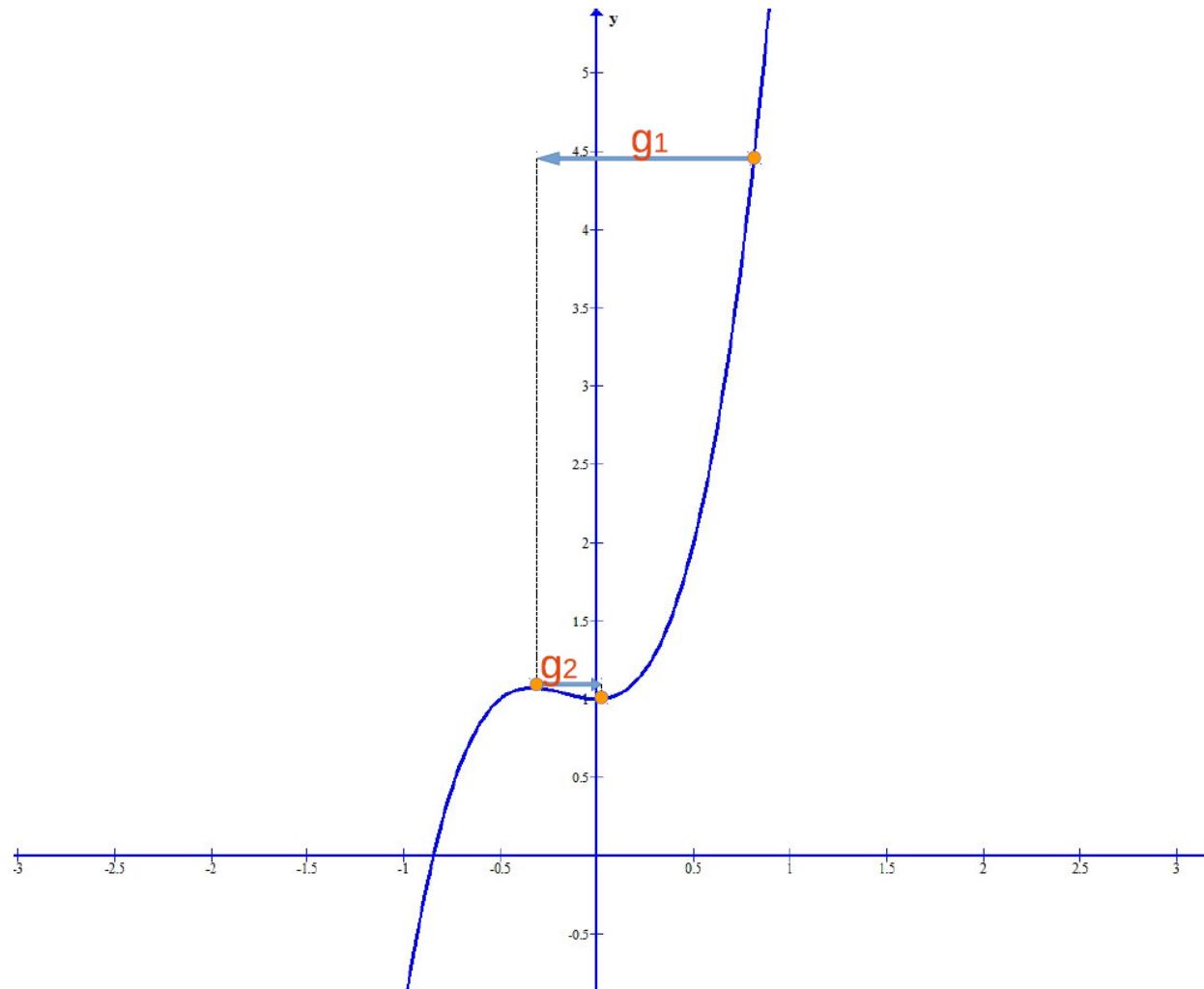
# Sigmoid



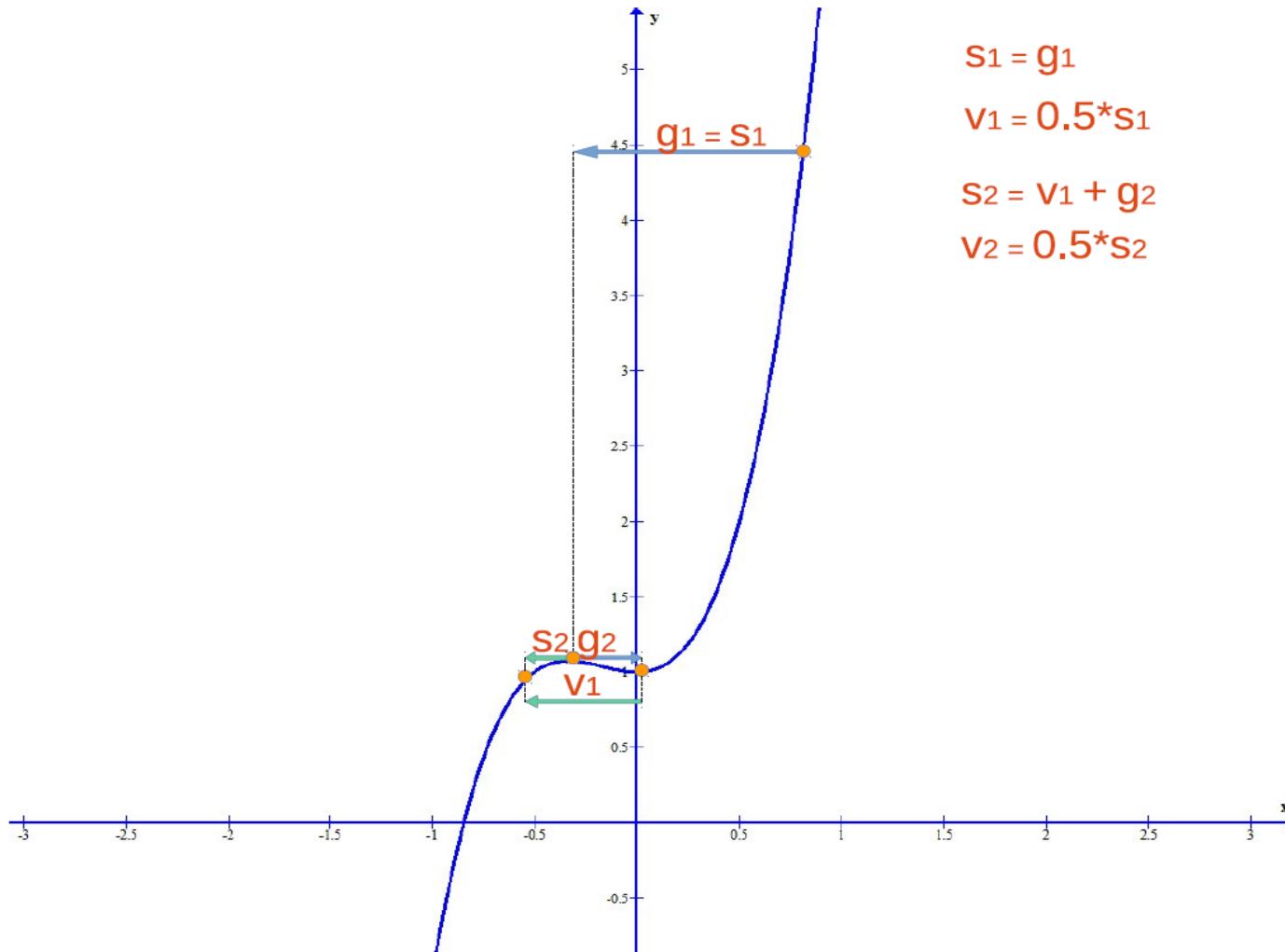
$$S(t) = \frac{1}{1 + e^{-t}}$$

- gradientni spust (minimiziramo funkciju napake)
  - $\Delta w_{ij} = \alpha \Delta_i x_j$        $\Delta_i = (y_i - S(in_i)) * S'(in_i)$
  - $\Delta b_i = \alpha \Delta_i$
  - $\alpha$  ... stopnja učenja
- vzratno širjenje napake
  - $\Delta_j = S'(in_i) * \sum_i w_{ij} \Delta_i$

# Gradientni spust



# Moment



- matematična operacija
- drseče okno (filter)
- vsota slikovnih elementov (pixlov) pomnoženih z utežmi filtra

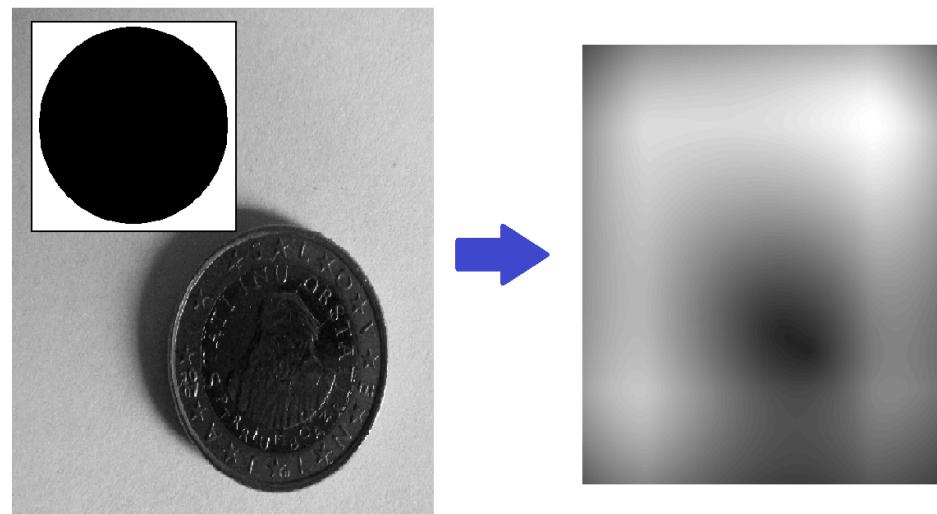
1	4	0
0	1	-1
3	-2	0

1	-1
0	2

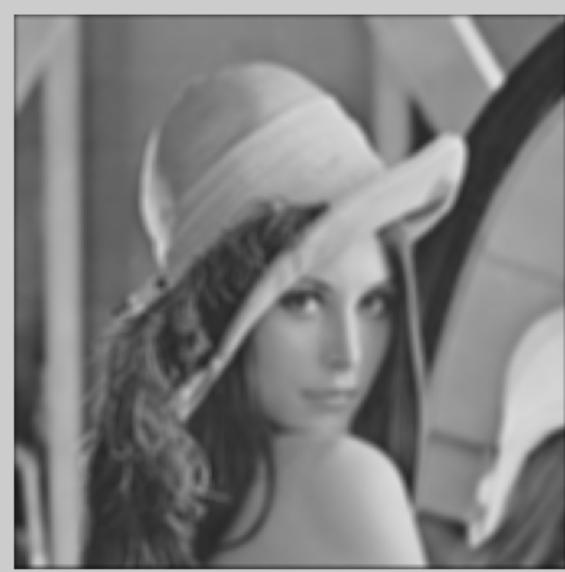
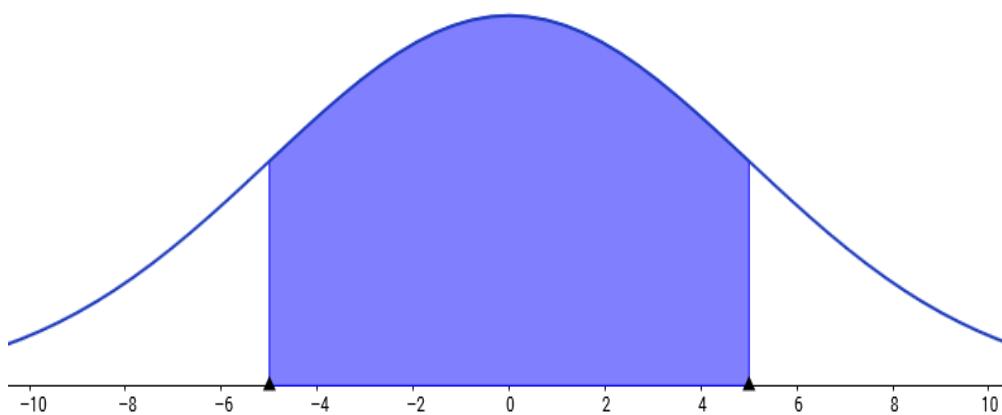


-1	2
-5	2

- detekcija objektov
- detektiranje robov
- transformacija slik



# Konvolucija



# Konvolucija

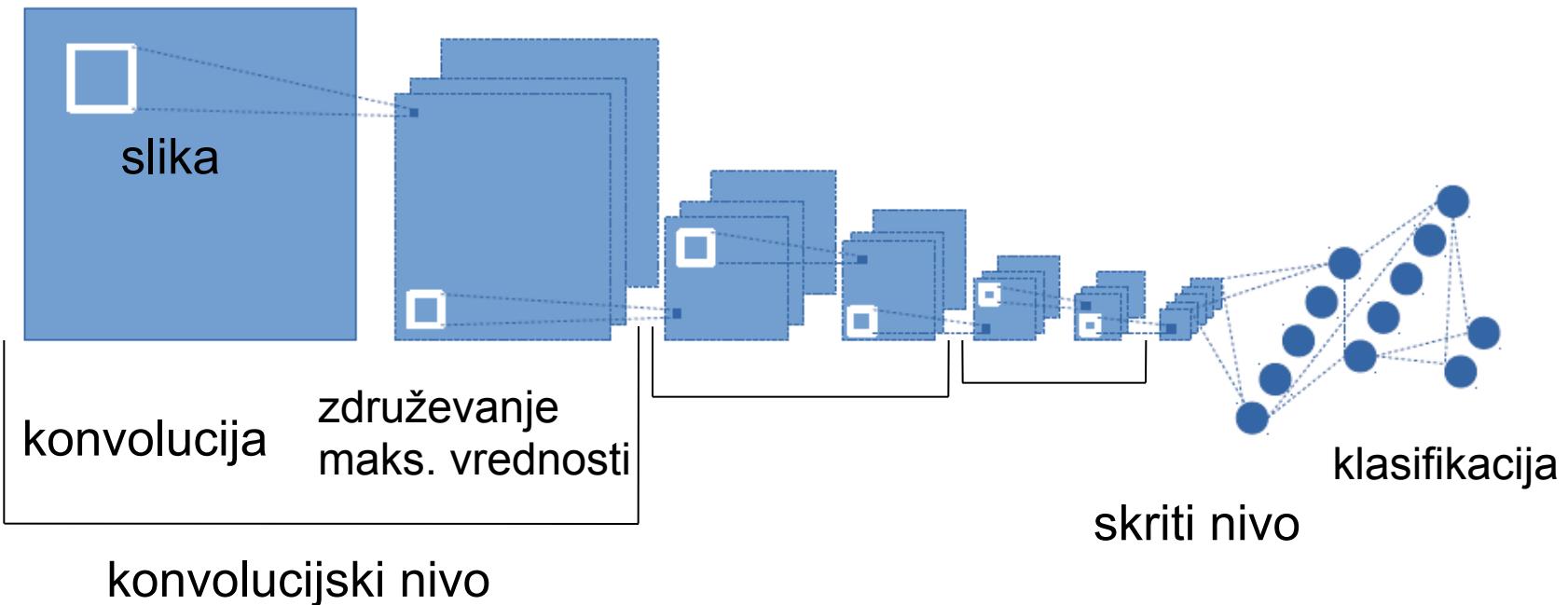
$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$



# Konvolucijske nevronske mreže

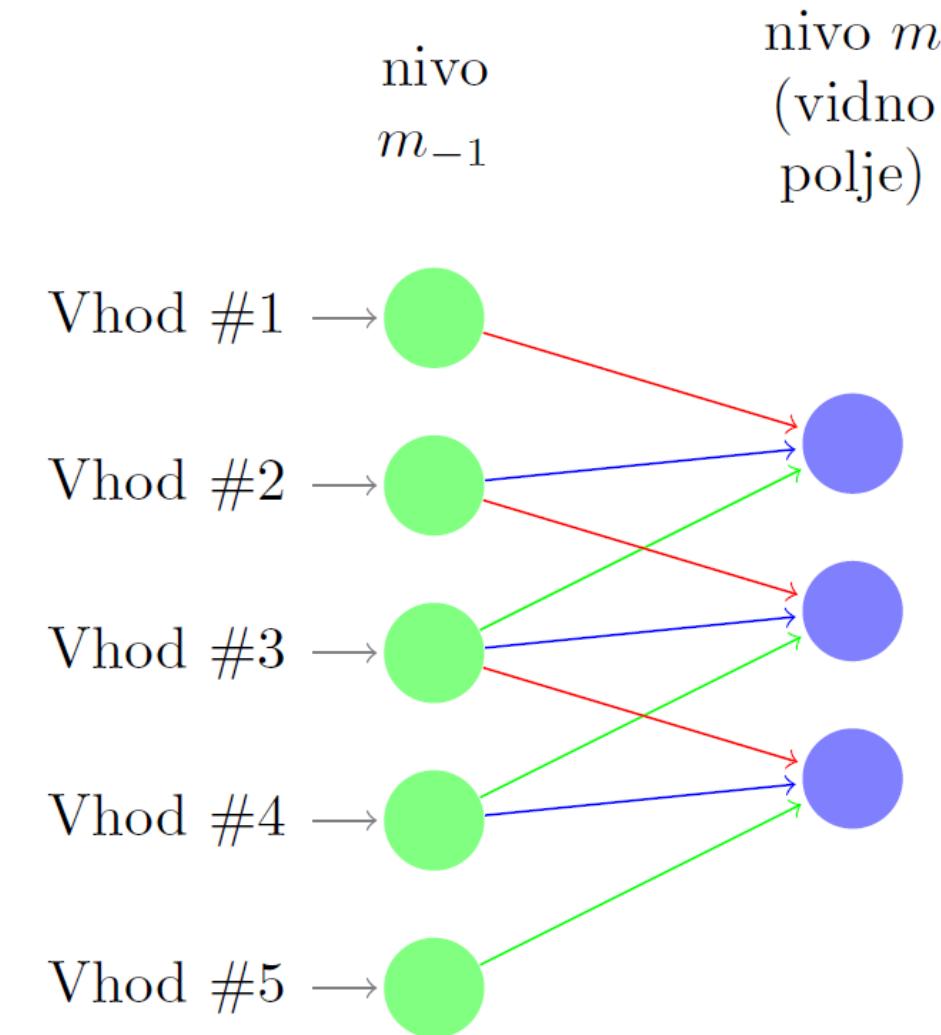
- namenjene vizualnemu zaznavanju
- nadgradnja preprostih nevronskih mrež
- globoke arhitekture
- veliko konvolucijskih nivojev
- prepoznavanje obrazov, ročno pisanih številk, razlikovanje med spoloma, ...

# Konvolucijske nevronske mreže



Pri združevanju maksimalnih vrednosti posredujemo samo največjo vrednost.

# Konvolucijske nevronske mreže



# Klasifikacija člankov

- članki pridobljeni z Dbpedie
- 10 razredov:
  - athlete, animal, musical work, village, artist  
film, infrastructure, building, company, natural  
place
- 30.000 učnih primerov
- 10.000 testnih primerov

- 150 zaporednih znakov iz vsakega članka
- lematizacija in krnjenje,
  - (NLTK korpusi, Lancaster stemmer)
- golo besedilo v ASCII kodiranju (bitno polje)
- vektorizacija z 1 do 45 kodiranjem, znaki
- a, b, c, d, e, f, g, h, i, j, k, l, m, n, o, p, q, r, s, t, u, v, w, x, y, z, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, '.', '!', '?', ',', '.', '-', "'", "(", ")"

# Obdelava člankov (primer)

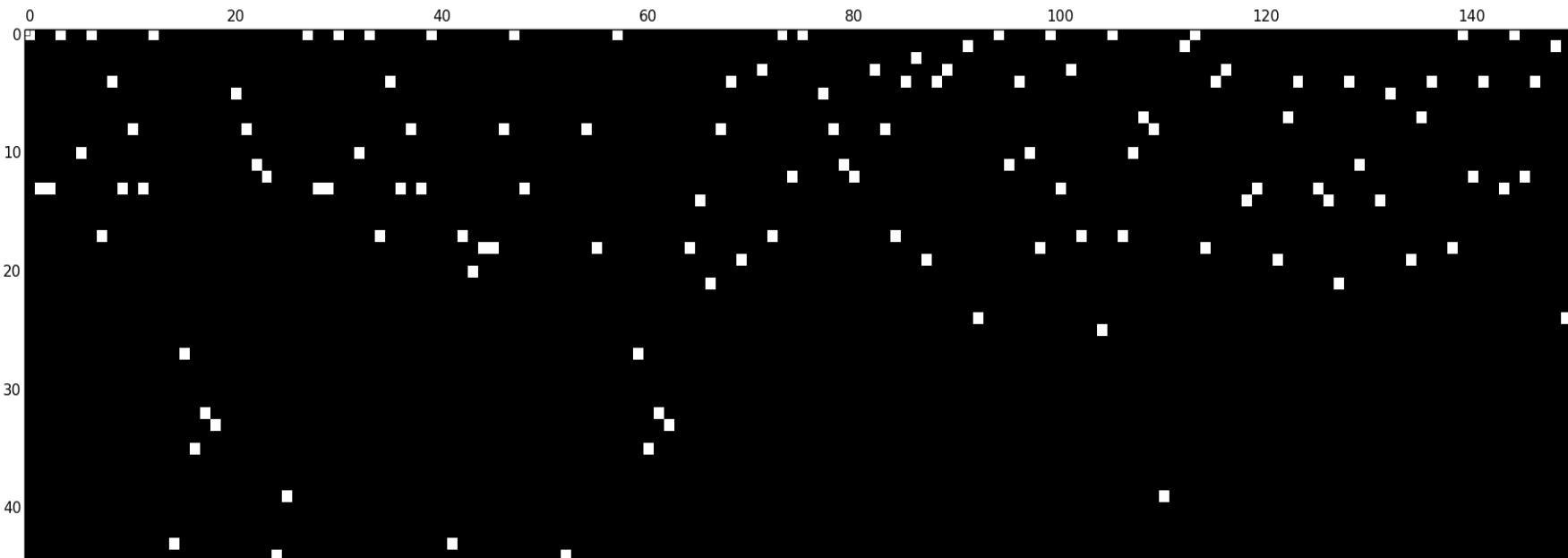
- Anna\_Karenina\_(1967\_film), Anna Karenina (Russian: Анна Каренина) is a 1967 Soviet drama film directed by Aleksandr Zarkhi, based on the novel of the same name by Leo Tolstoy. It was listed to compete at the 1968 Cannes Film Festival, but the festival was cancelled due to the events of May 1968 in France.

# Obdelava člankov (primer)

- anna\_karenina\_(1967\_film) , ann karenin ( russ : ) is a 1967 soviet dram film direct by aleksandr zarkh , bas on the novel of the sam nam by leo tolstoy . it was list to compet at the 1968 can film fest , but the fest was cancel due to the ev of may 1968 in frant .

# Obdelava člankov

- podobnost z Brajevo pisavo

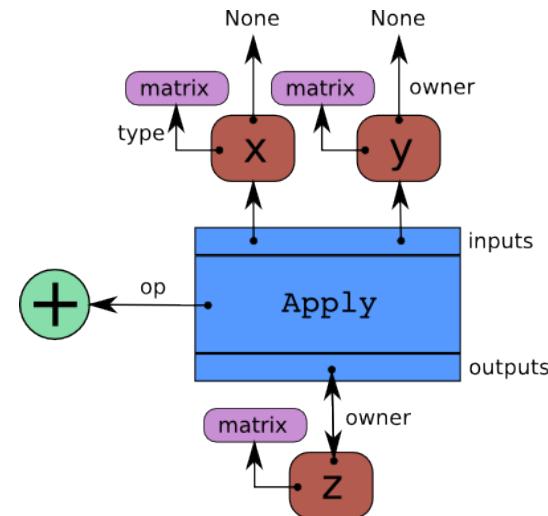


# Obdelava člankov

- primer: 1 - 7 kodiranje, "c, e, i, j, n, o, s"

s	o	n	c	e	s	i	j	e
o	o	o	1	o	o	o	o	o
o	o	o	o	1	o	o	o	1
o	o	o	o	o	o	1	o	o
o	o	o	o	o	o	o	1	o
o	o	1	o	o	o	o	o	o
o	1	o	o	o	o	o	o	o
1	o	o	o	o	o	1	o	o

- programski jezik python
- knjižnjica Theano
  - optimizacija simboličnih izrazov
  - simbolično odvajanje

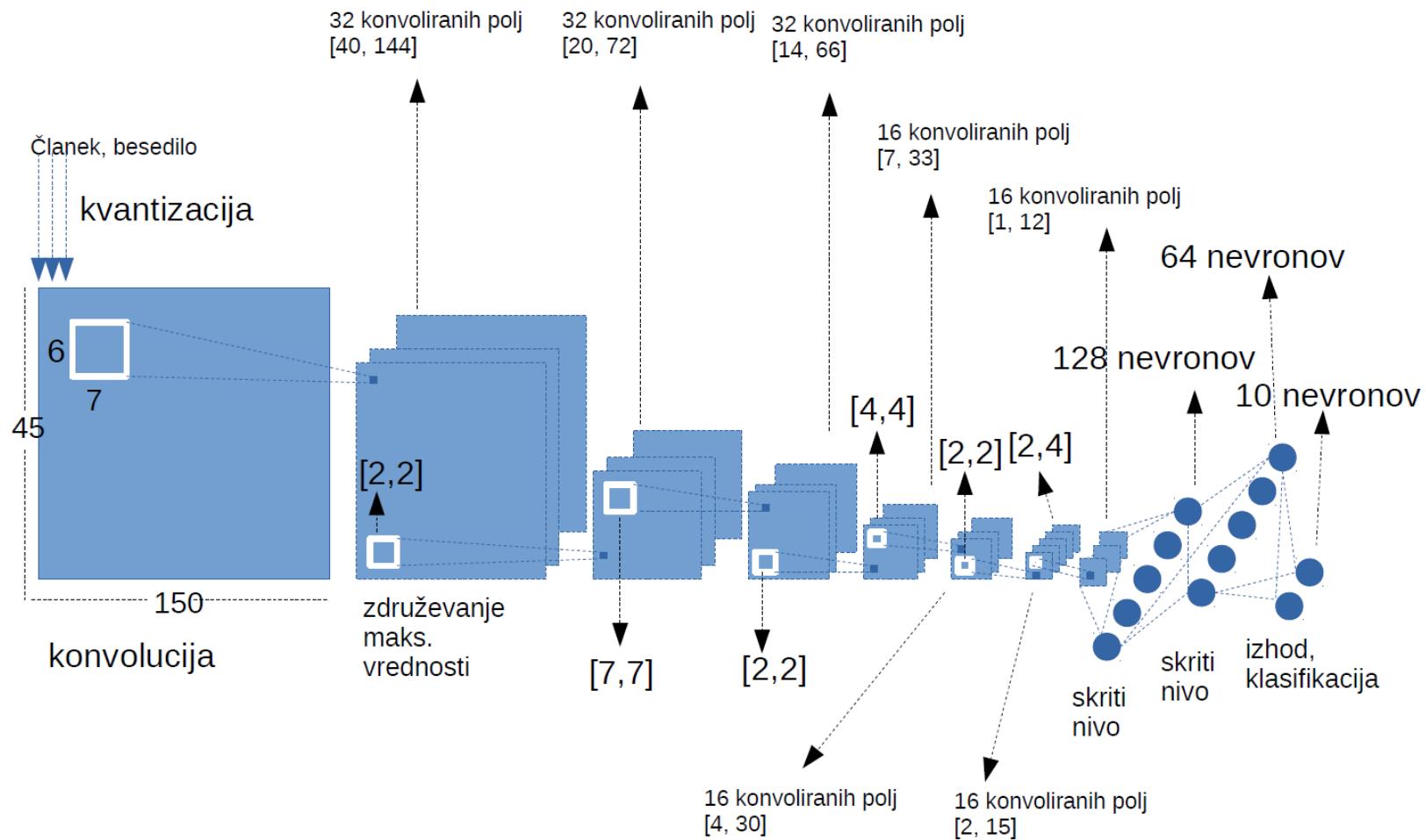


[http://deeplearning.net/software/theano/tutorial/symbolic\\_graphs.html](http://deeplearning.net/software/theano/tutorial/symbolic_graphs.html)

# Razvoj ideje in implementacija

- Kako globoko arhitekturo potrebujemo?
- Kakšne so primerne dimenzijske razmerje filterov?
- model za kvantizirane podatke in golo besedilo
  - 4 konvolucijski nivoji, 2 skrita nivoja, izhodni nivo
  - filtri za golo besedilo so enodimensionalni

# Model



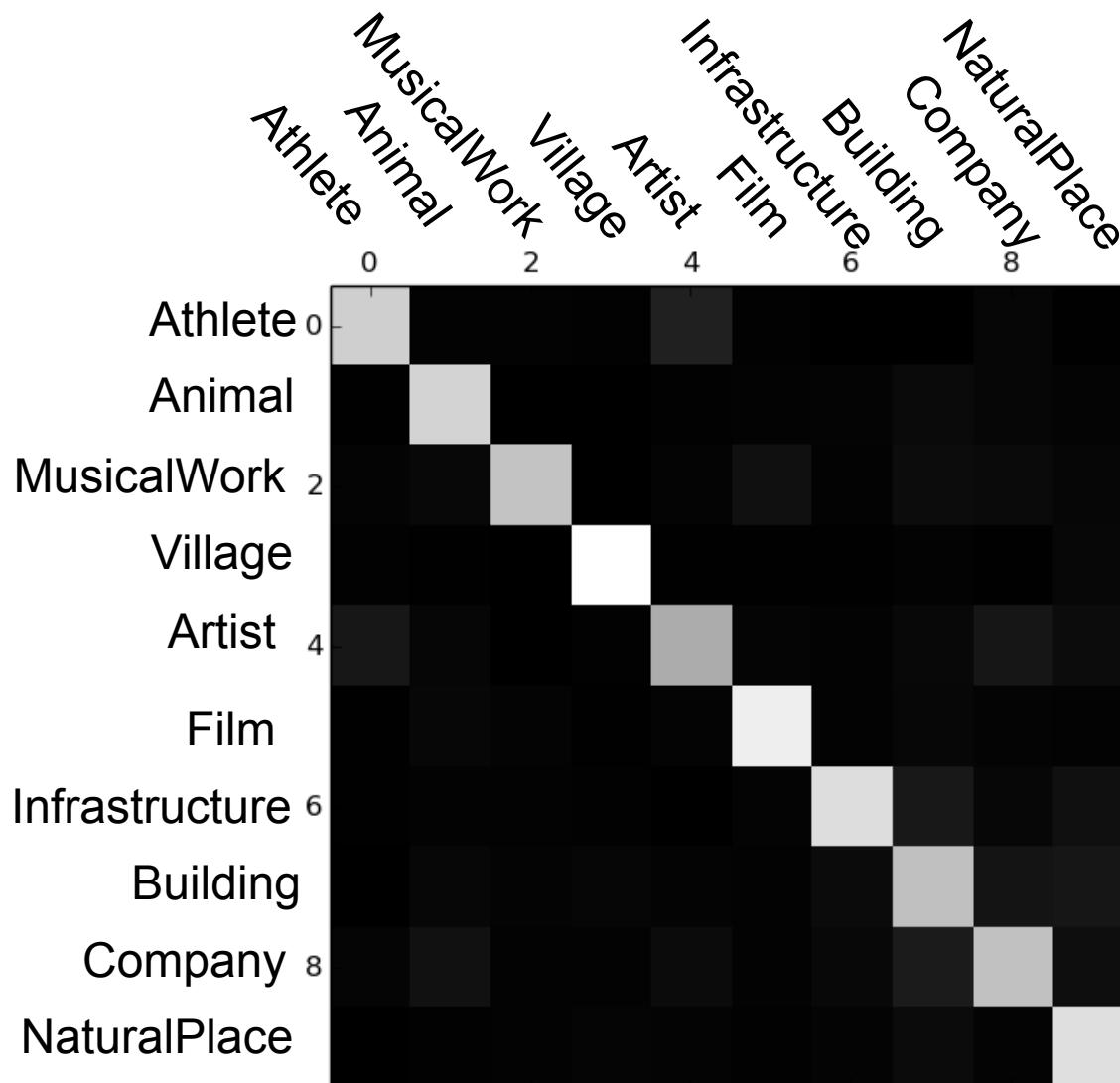
# Rezultati - kategorizacija člankov

za stopnja učenja = 0.13 in stopnja momenta 0.04

klasifikacijska točnost (%)	krnjenje	lematizacija
bitno polje (ASCII koda)	20.79%	46.33%
vektorizacija	82.90%	80.90%

- dvodimensionalna konvolucija prekaša enodimensionalno

# Matrika zmot



- velik potencial za jezikovne probleme
- primerne za malo obdelane podatke
- poiskati ustrezeno parametrizacijo
- časovno potratne
- računsko zahtevne
- prilagoditve arhitekture
- več raziskav
- še globlje konvolucijske nevronske mreže

- klasifikacija člankov
- „sentiment analysis“
- učenje jezikovnih pravil (vejice)
- razčlenjevanje stavkov
- ...



Vprašanja?