

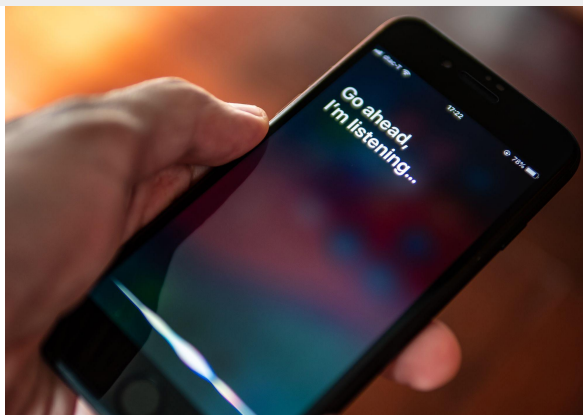
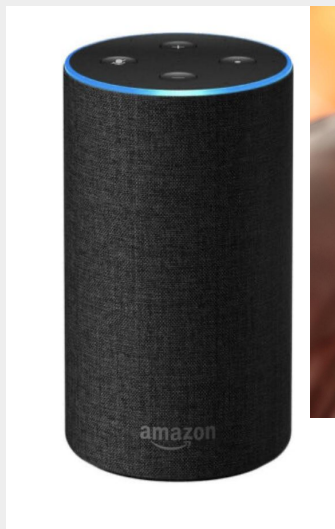
Može li AI da prepozna emocije?  
Razumevanje govora

Kristina Nikolić

Matematička gimnazija

16. 05. 2023.

1. Uvod u obradu govora
2. Metode dubokog učenja
3. Prepoznavanje emocija

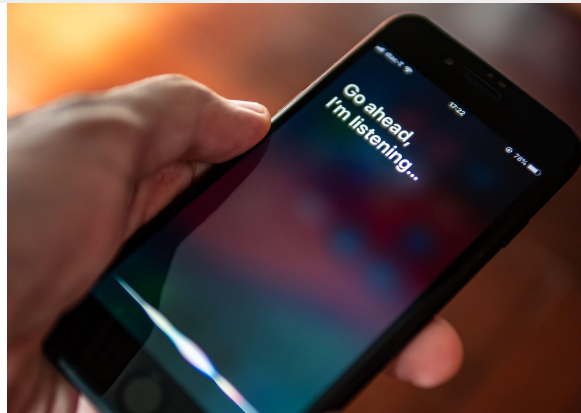
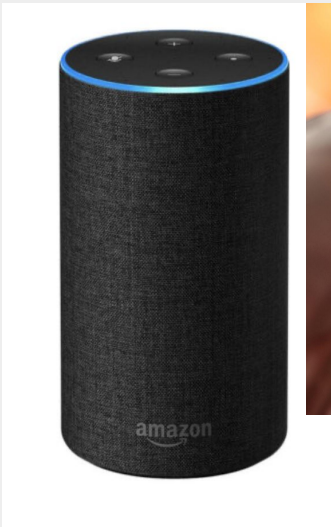


Govorni asistenti

# Primene razumevanja govornog signala



Hands free sistem u automobilima

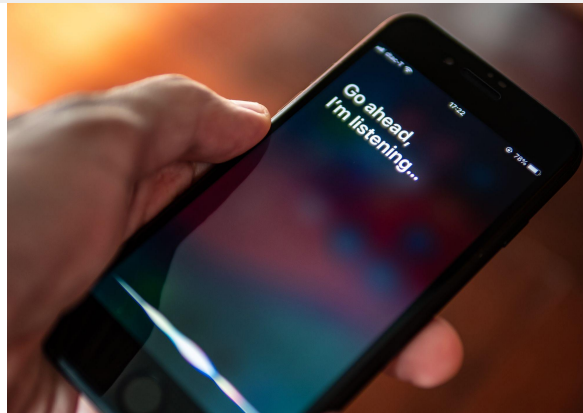
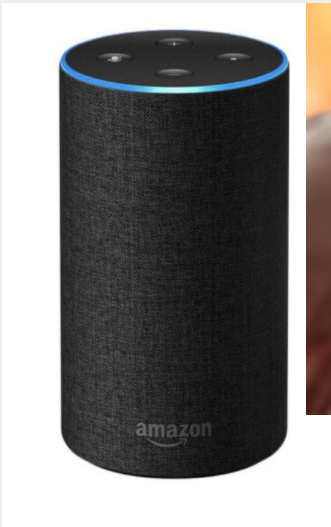


Govorni asistenti

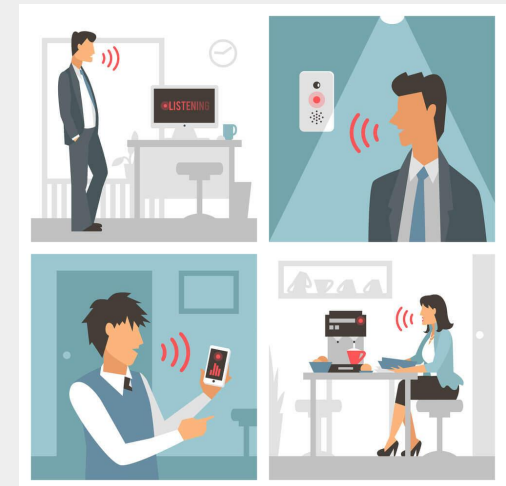
# Primene razumevanja govornog signala



Hands free sistem u automobilima



Govorni asistenti

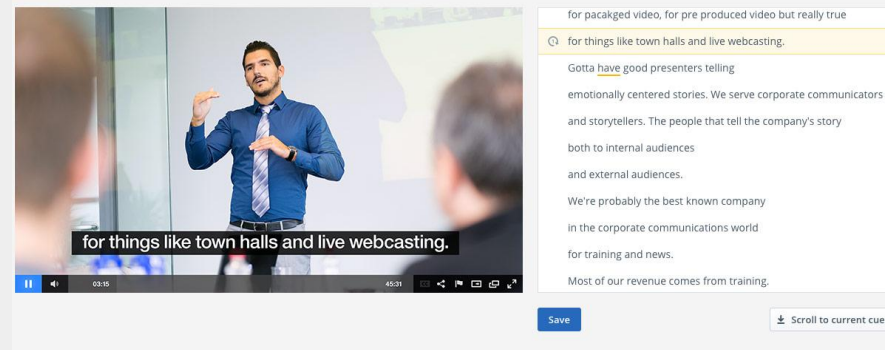


Smart home

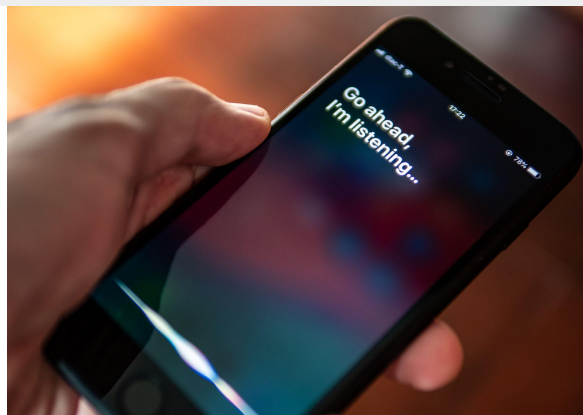
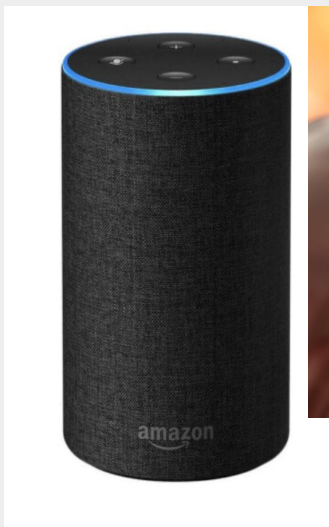
# Primene razumevanja govornog signala



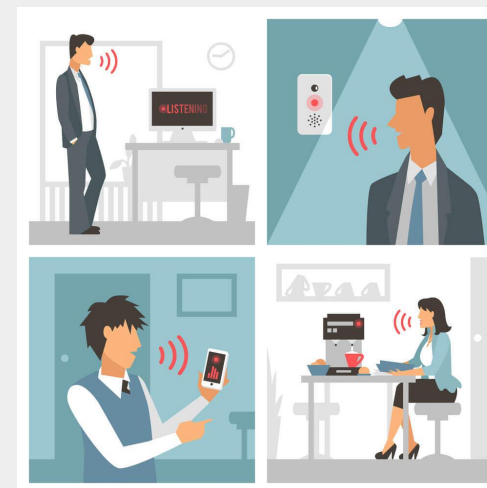
Hands free sistem u automobilima



Automatski transkript

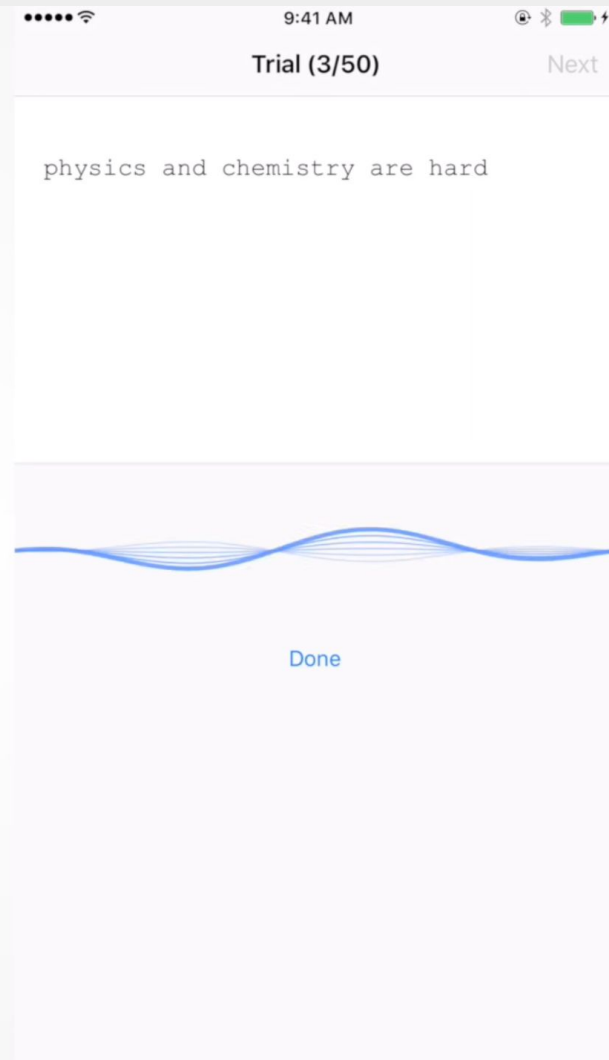
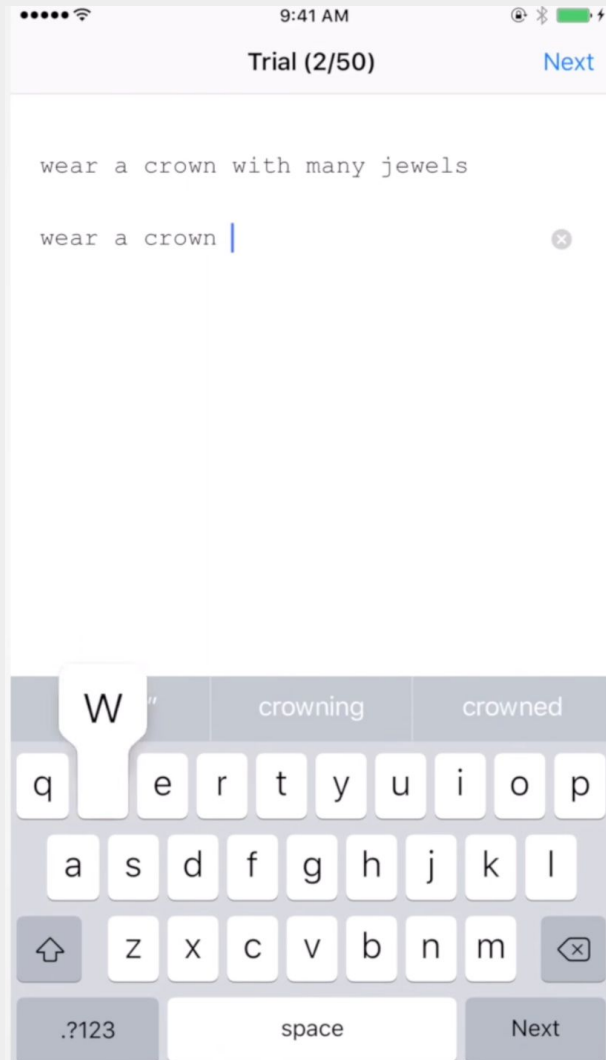


Govorni asistenti



Smart home

# Stanford: „Govorne komande 3x brže od kucanja”





- Govor u tekst
- Detekcija ključnih reči (word spotting)
- Identifikacija govornika (Verifikacija)



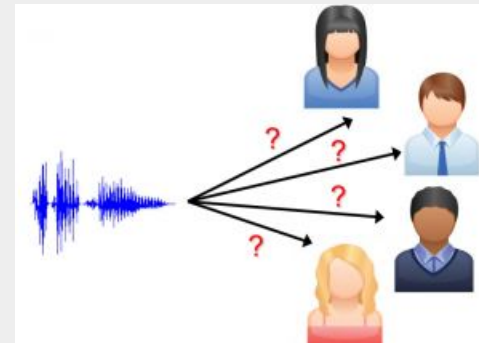


- Govor u tekst
- Detekcija ključnih reči (word spotting)
- Identifikacija govornika (Verifikacija)





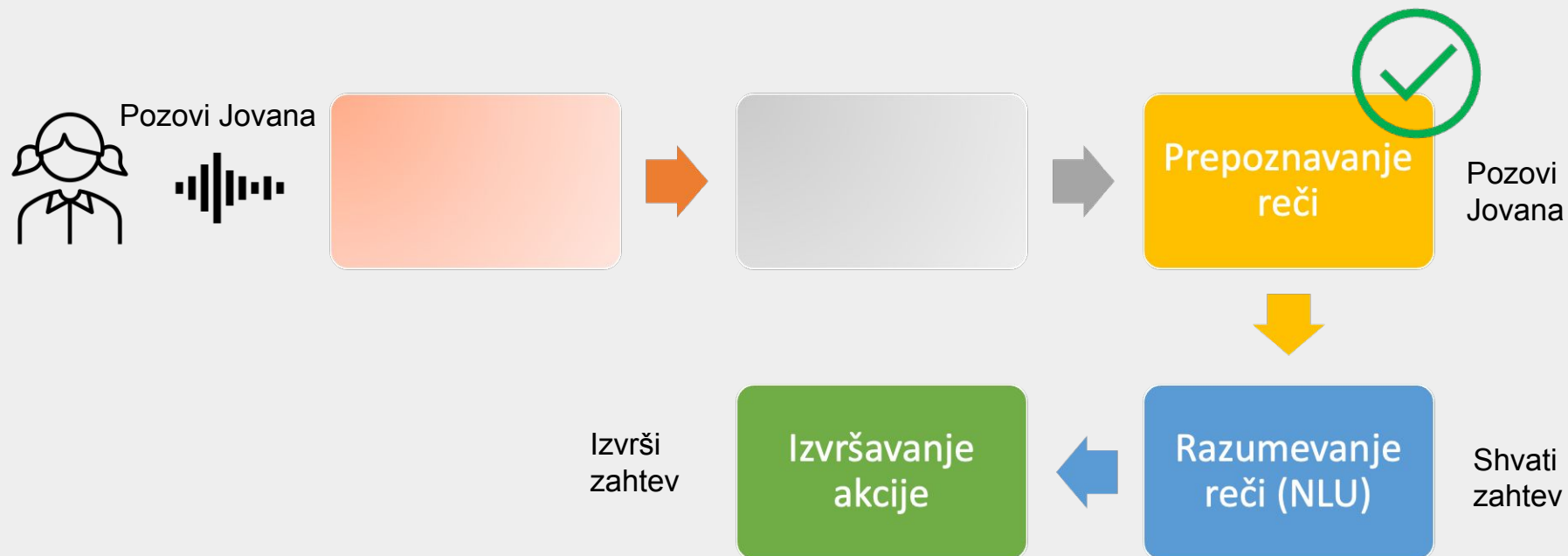
- Govor u tekst
- Detekcija ključnih reči (word spotting)
- Identifikacija govornika (Verifikacija)



# Kako rade govorni asistenti?



# Kako rade govorni asistenti?



# Kako rade govorni asistenti?



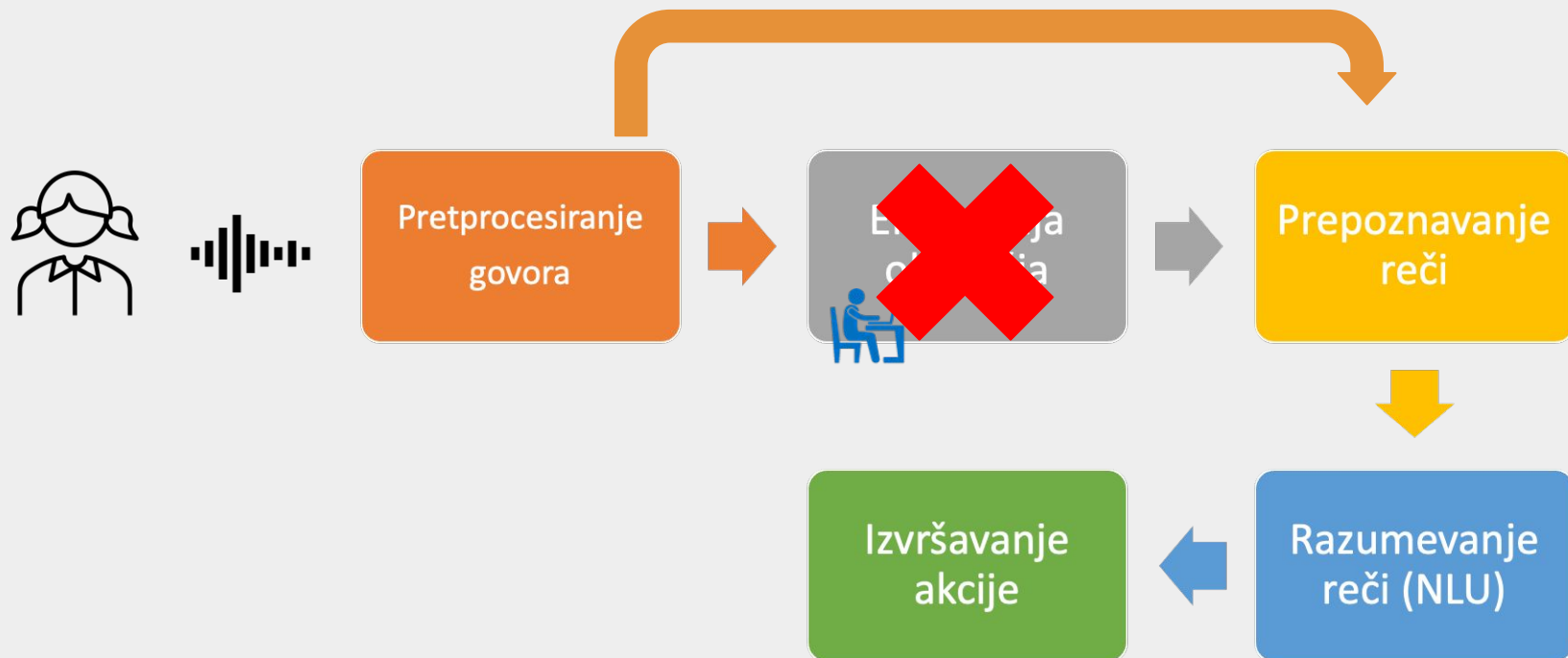
- Preporcesiranje
  - Uklanjanje šuma, augmentacija...
  - Spektogram, Mel spektogram...

# Kako rade govorni asistenti?

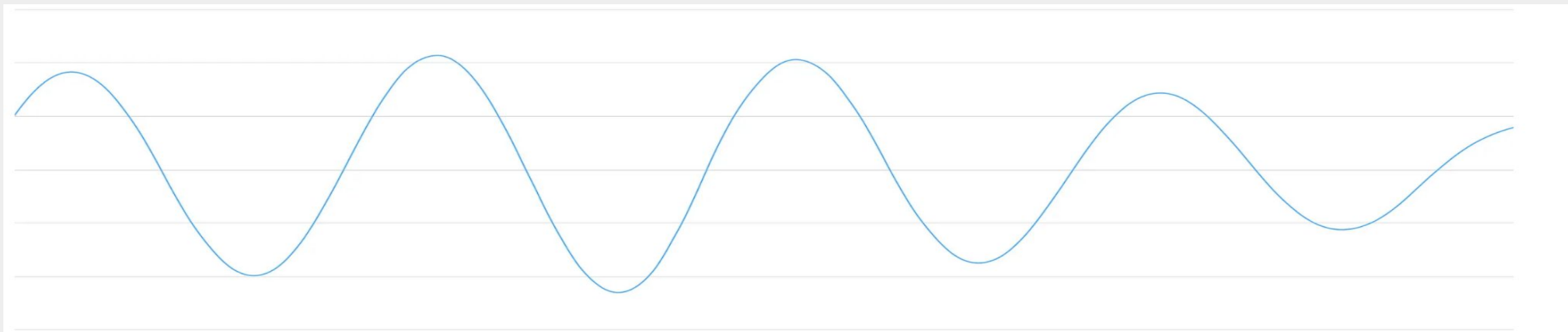


- Ekstrakcija obeležja
  - Izdvajamo domensko znanje koje smatramo korisnim
    - Standardna devijacija, snaga signala...
  - Metode dubokog učenja teže da izbegnu ovaj korak.

# Kako rade govorni asistenti?



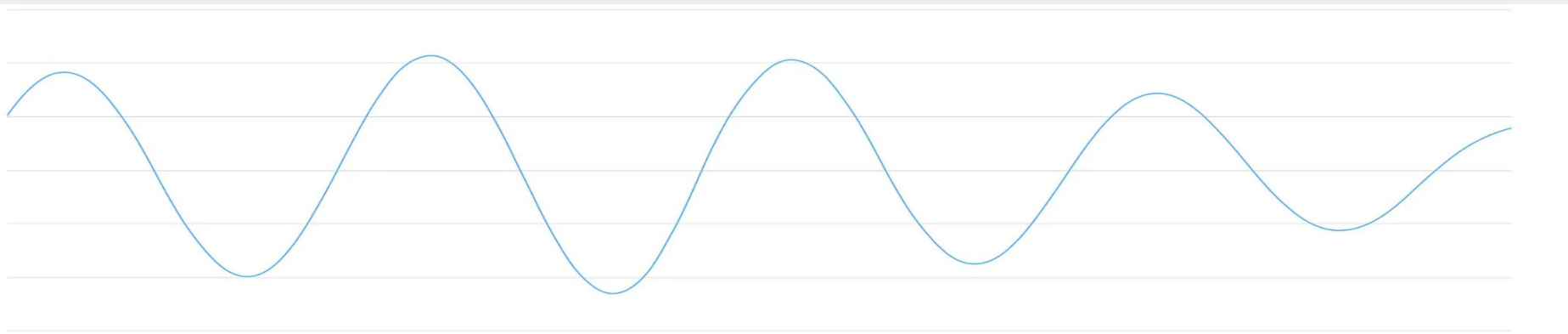
- Ekstrakcija obeležja
  - Izdvajamo domensko znanje koje smatramo korisnim
    - Standardna devijacija, snaga signala...
  - Metode dubokog učenja teže da izbegnu ovaj korak



- Želimo da predstavimo signal diskretno
- Uobičajna frekvencije odabiranja je 16kHz



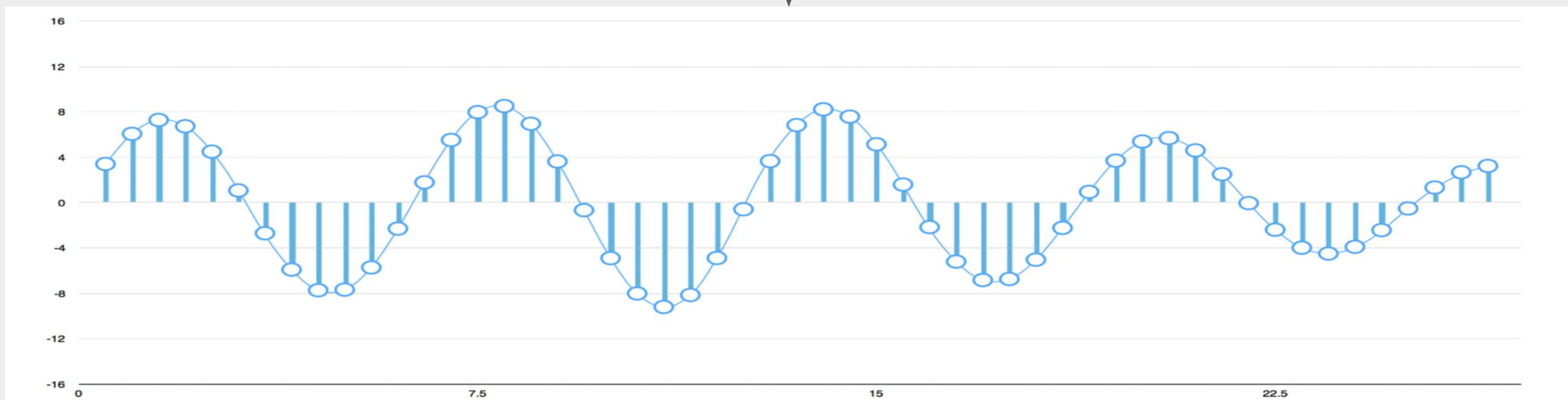
# Uzorkovanje govornog signala



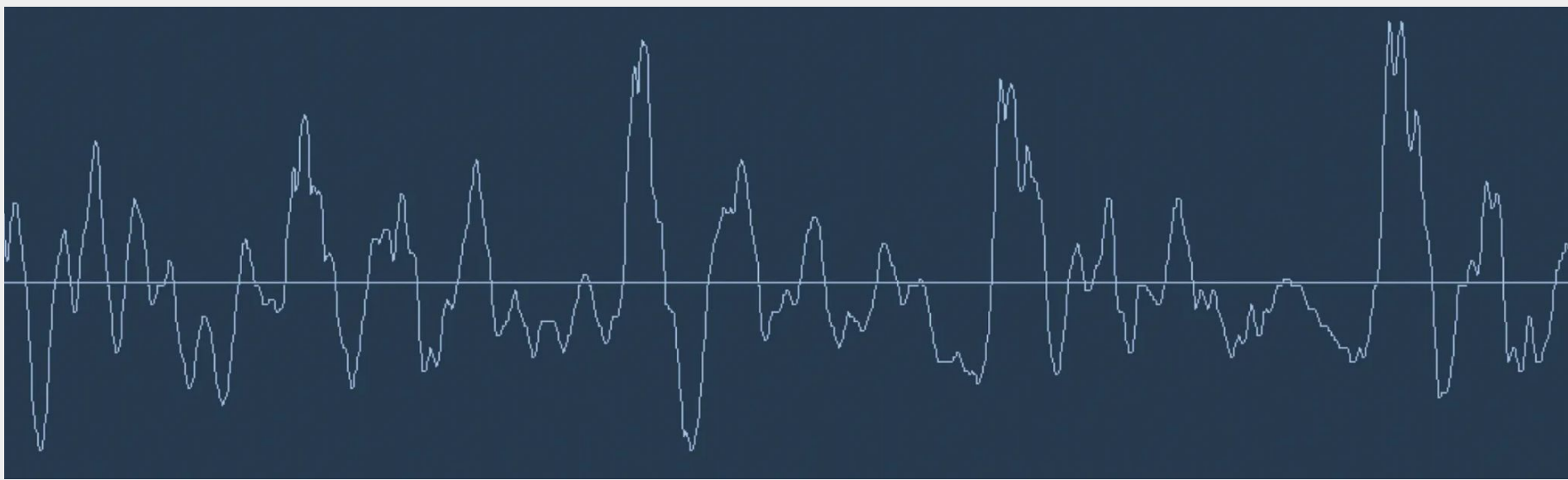
- Želimo da predstavimo signal diskretno



- Uobičajna frekvencije odabiranja je 16kHz



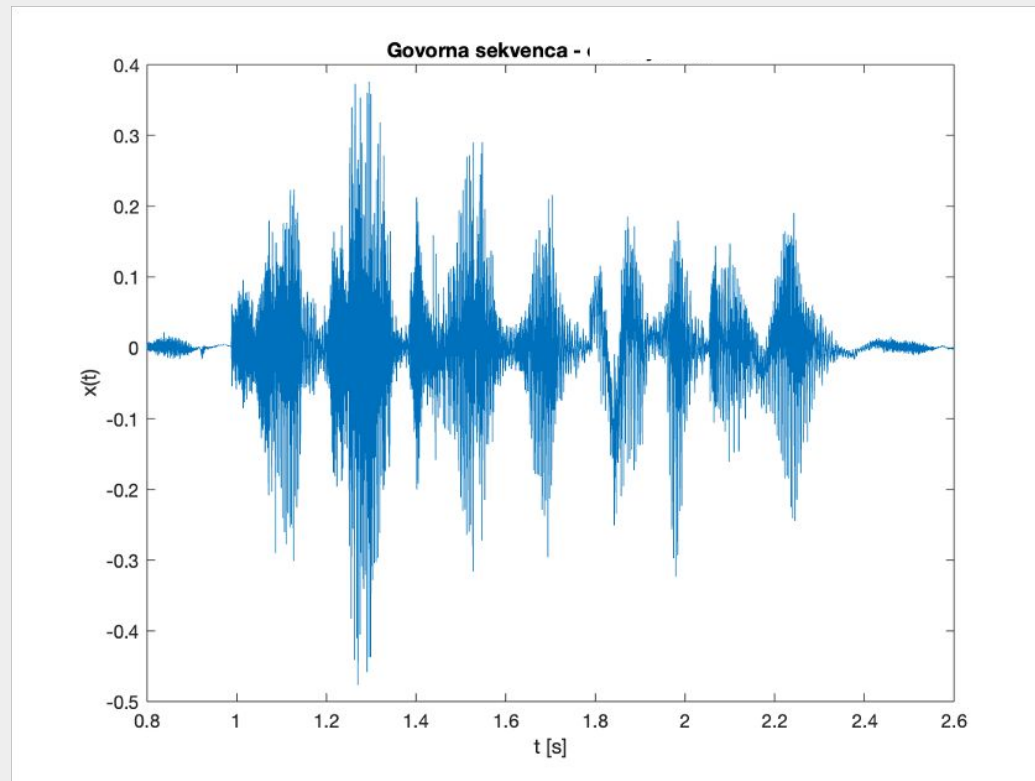
- Govorni signal reči: „Hello”:



- Pretstavljamo sekvencom brojeva (vrednosti signala u trenutku odabiranja)

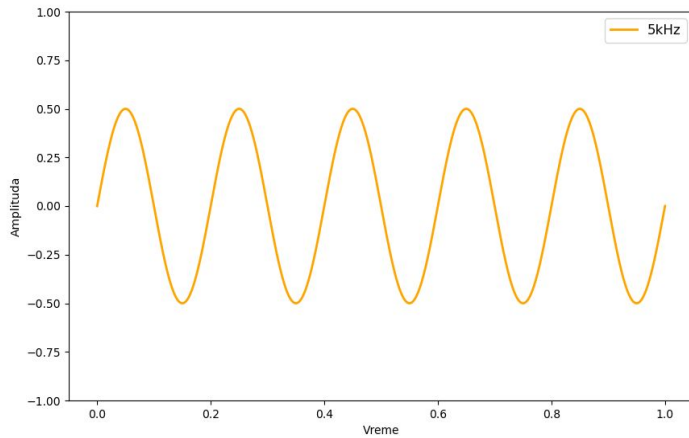
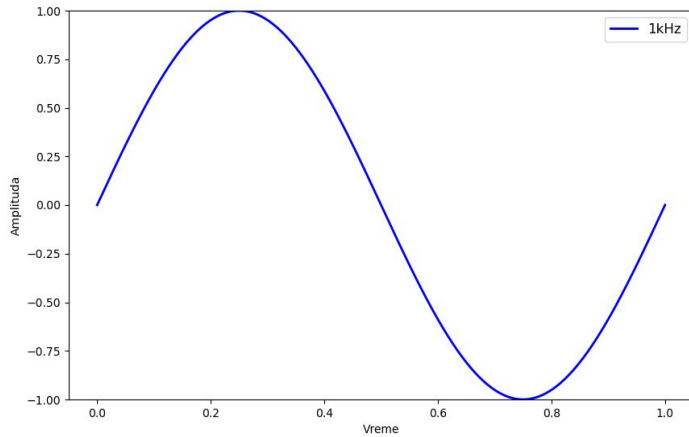
```
[-1274, -1252, -1160, -986, -792, -692, -614, -429, -286, -134, -57, -41, -169, -456, -450, -541, -761, -1067, -1231, -1047, -952, -645, -489, -448, -397, -212, 193, 114, -17, -110, 128, 261, 198, 390, 461, 772, 948, 1451, 1974, 2624, 3793, 4968, 5939, 6057, 6581, 7302, 7640, 7223, 6119, 5461, 4820, 4353, 3611, 2740, 2004, 1349, 1178, 1085, 901, 301, -262, -499, -488, -707, -1406, -1997, -2377, -2494, -2605, -2675, -2627, -2500, -2148, -1648, -970, -364, 13, 260, 494, 788, 1011, 938, 717, 507, 323, 324, 325, 350, 103, -113, 64, 176, 93, -249, -461, -606, -909, -1159, -1307, -1544]
```

- Primer govornog signala u **vremenskom domenu**:

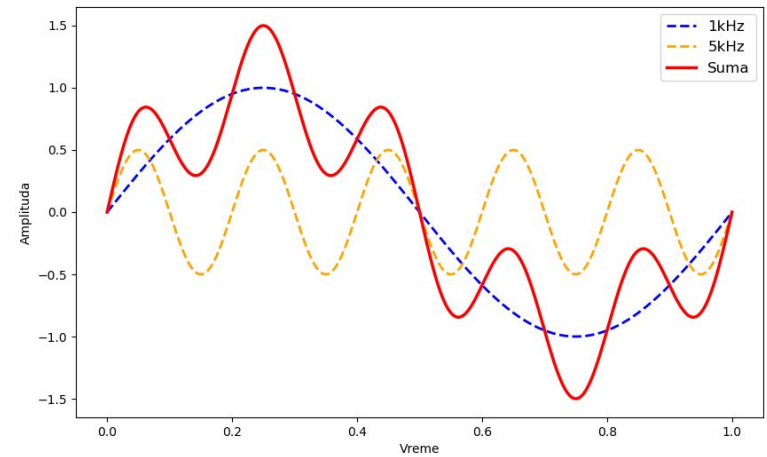
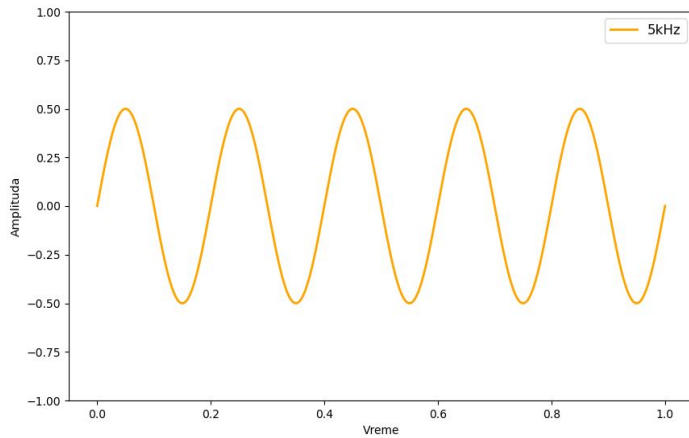
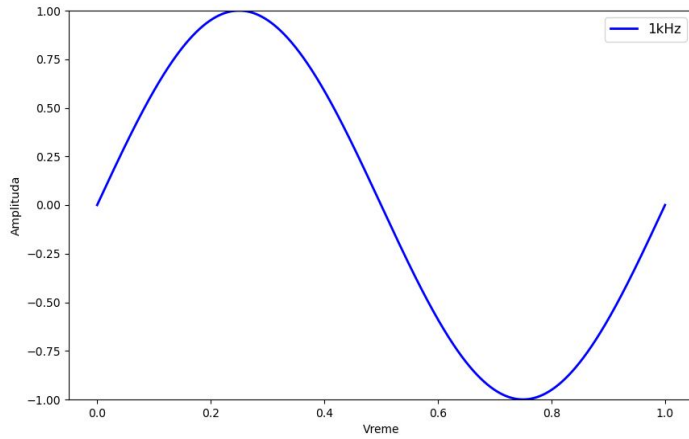


- Želimo da rastavimo ovaj signal na frekvencije koje ga čine, tj. da ga pretstavimo u **frekvencijskom domenu**.

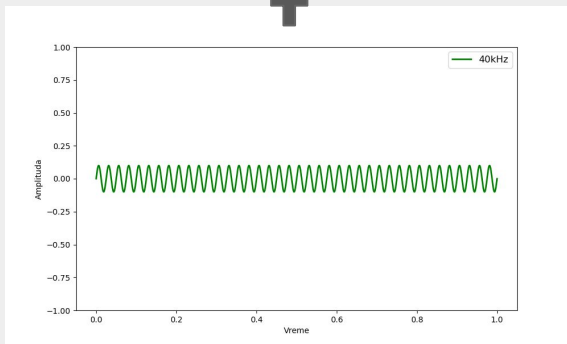
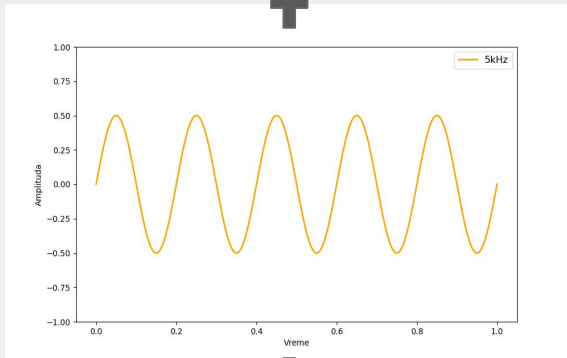
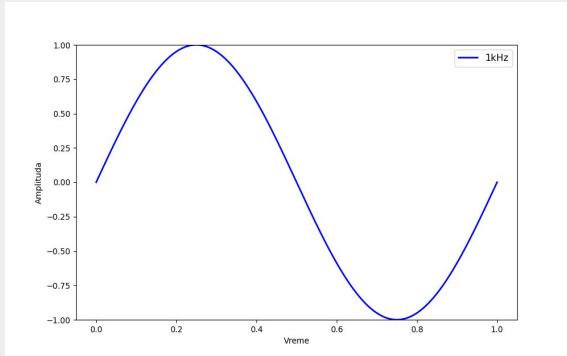
# Periodični signali



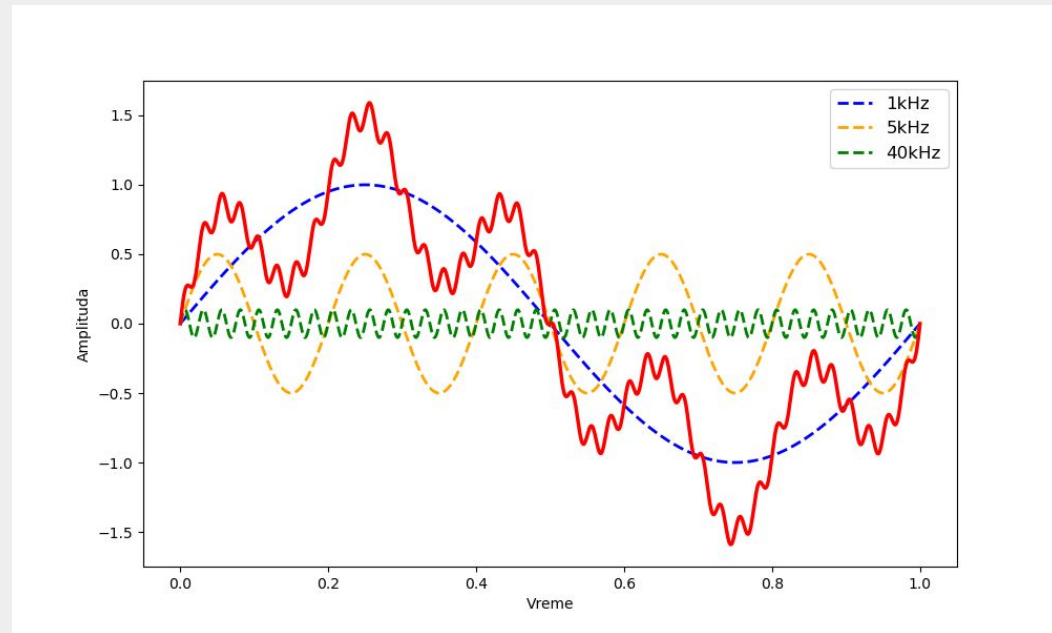
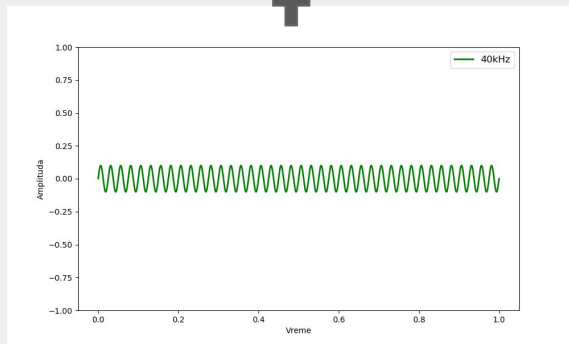
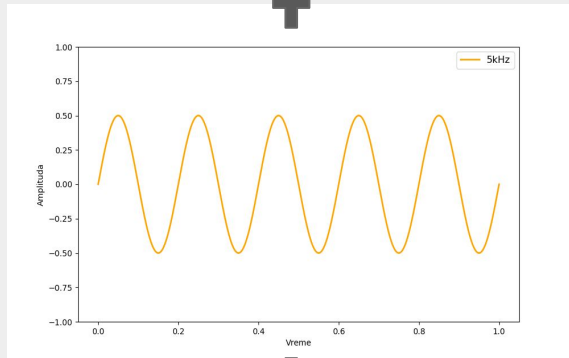
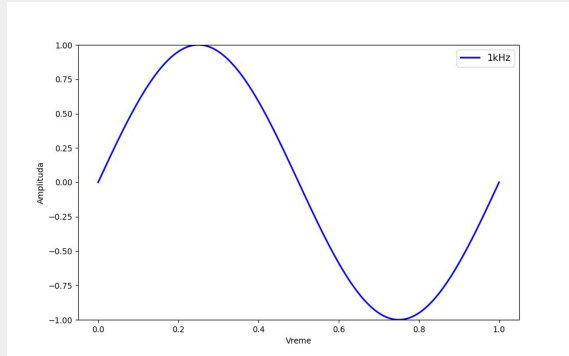
# Periodični signali



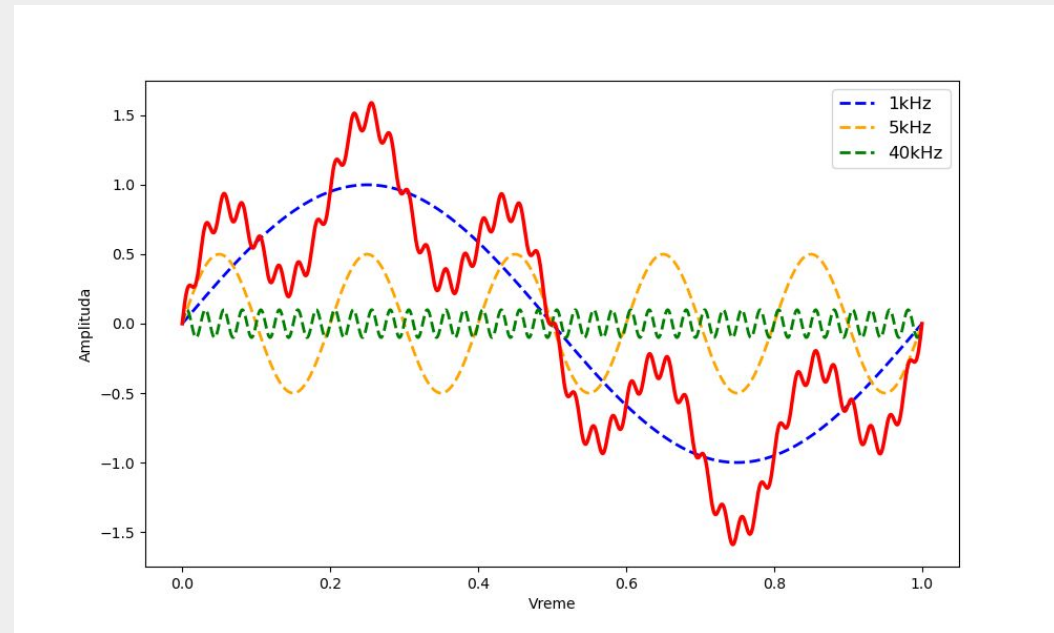
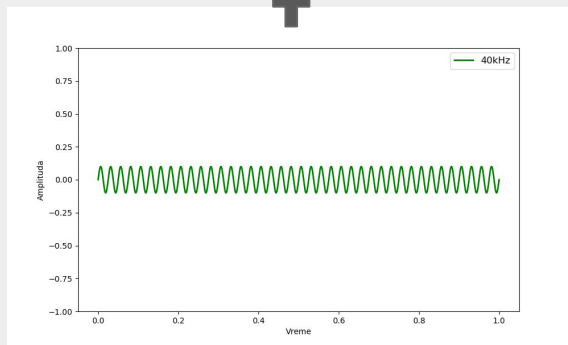
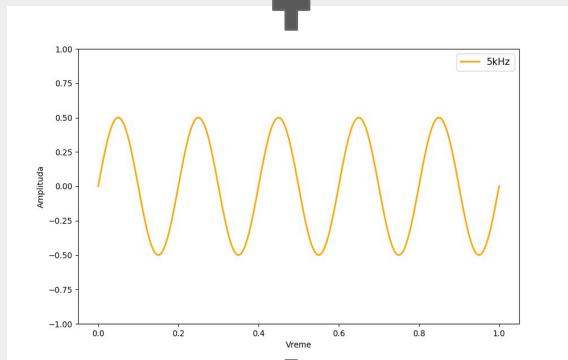
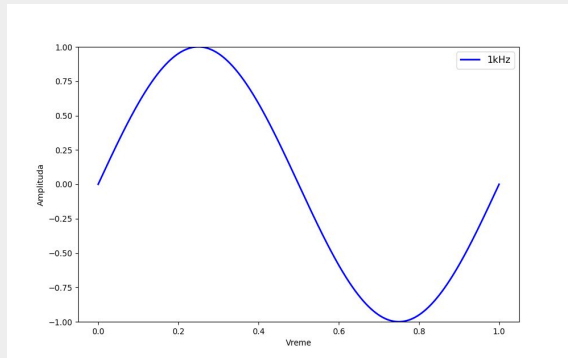
# Periodični signali



# Periodični signali

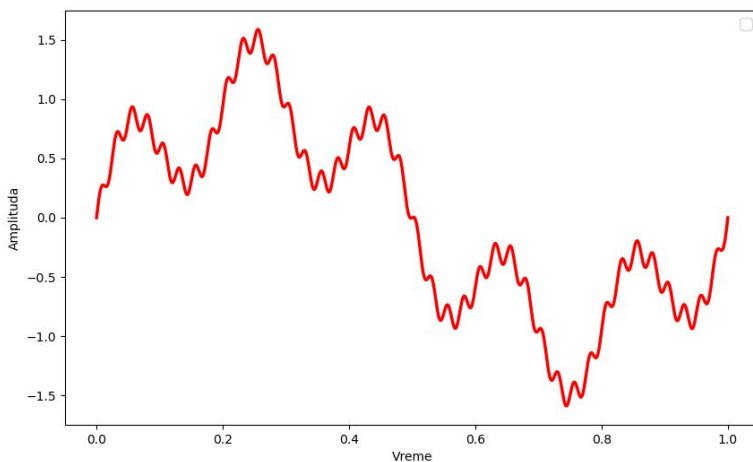


# Periodični signali

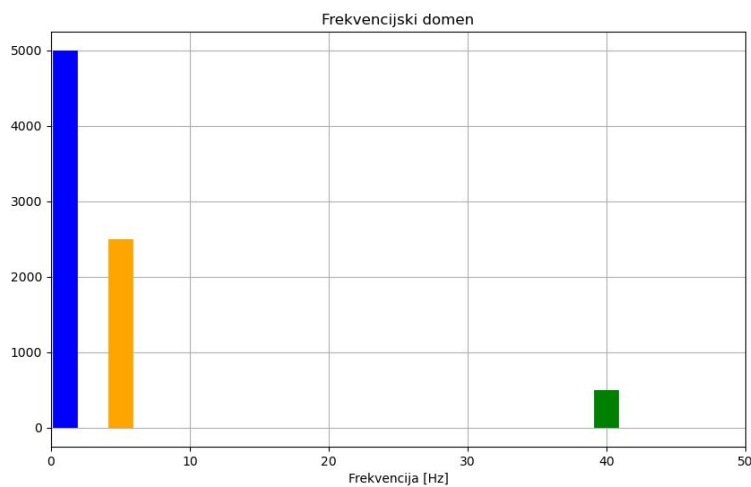


Matematički alat koji nam omogućava da periodične signale rastavimo na sumu prostih signala (sin i cos) naziva se **Furijeov red**.



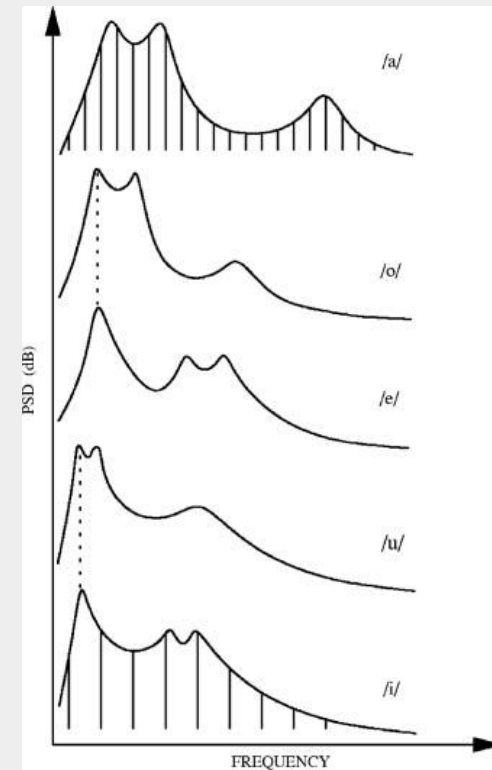
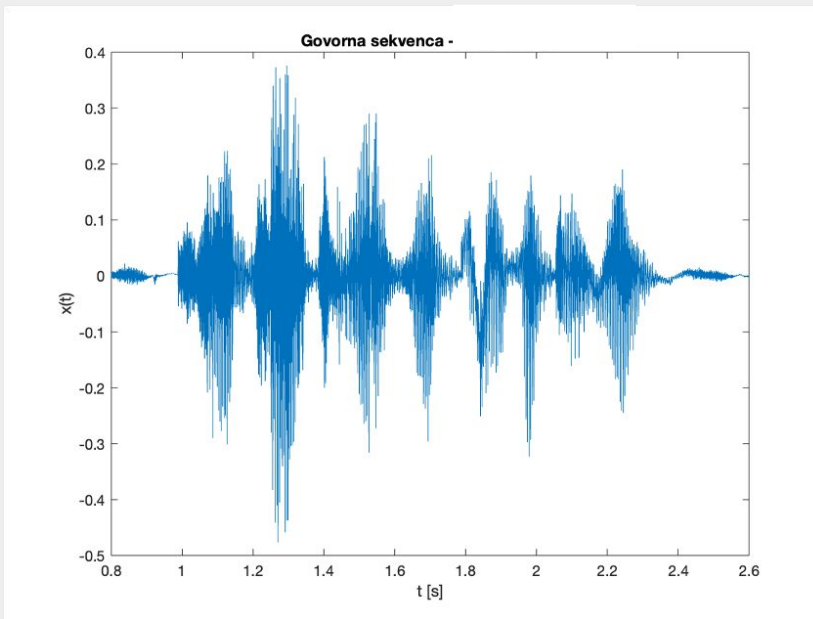


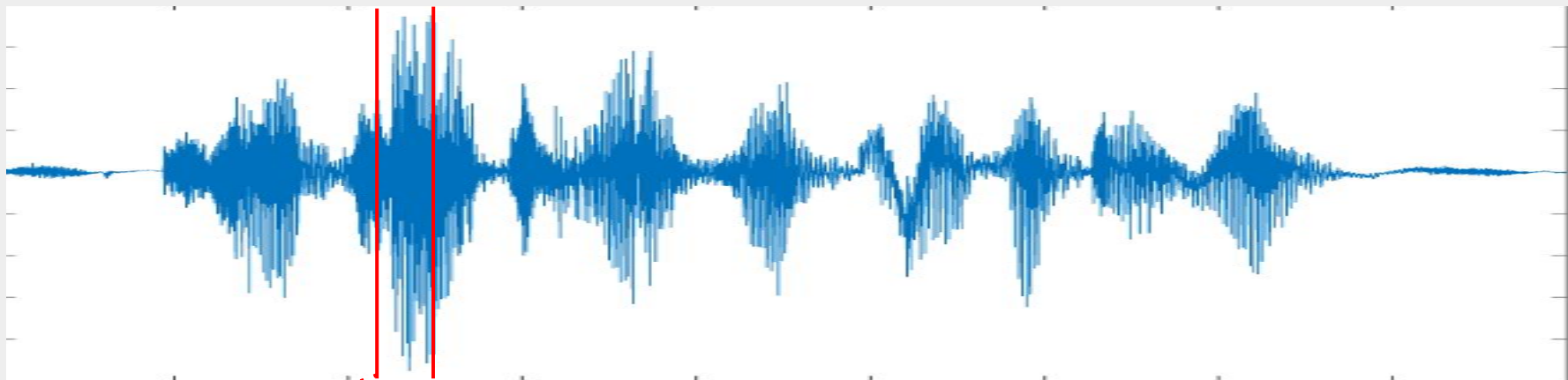
**Vremenskom domen**



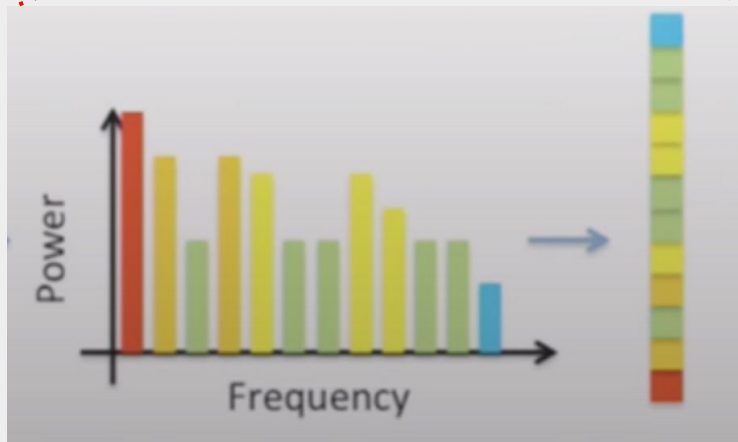
**Frekvencijski domen**

- Signali u prirodi nisu periodični, stoga njihov spektralni domen nije diskretan, tj. snaga signala je zastupljena na svim frekvencijama.
- Za spektralnu analizu neperiodičnih signala koristi se **Furijeova transformacija**.



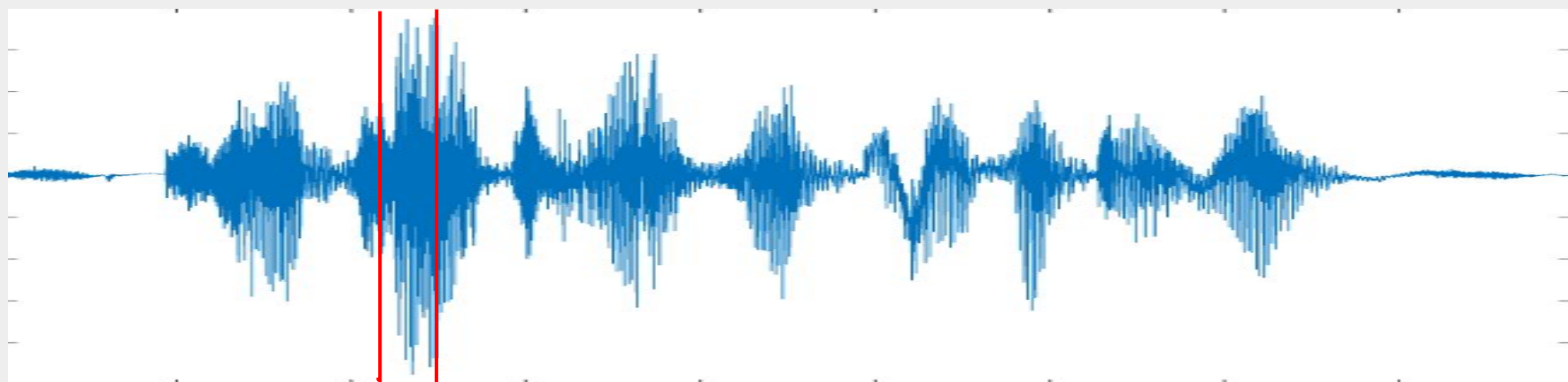


20ms

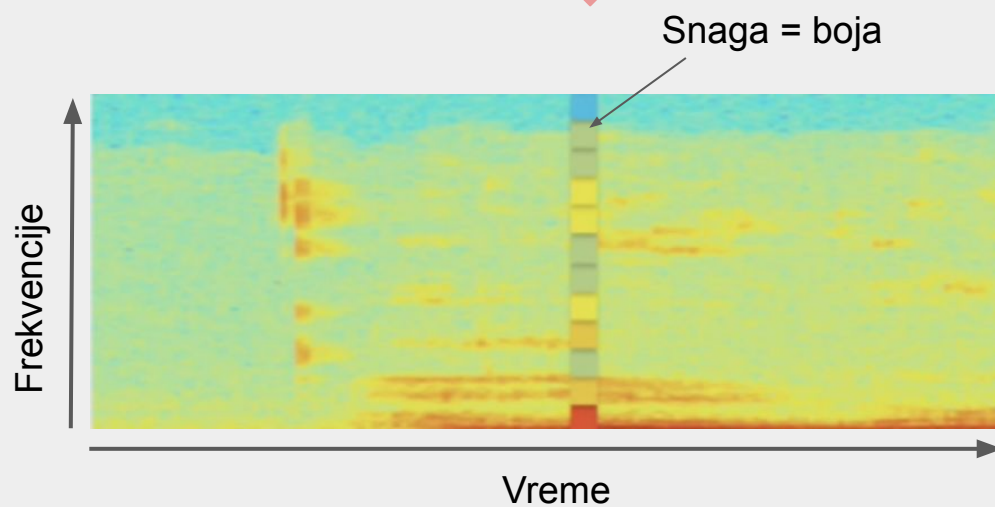
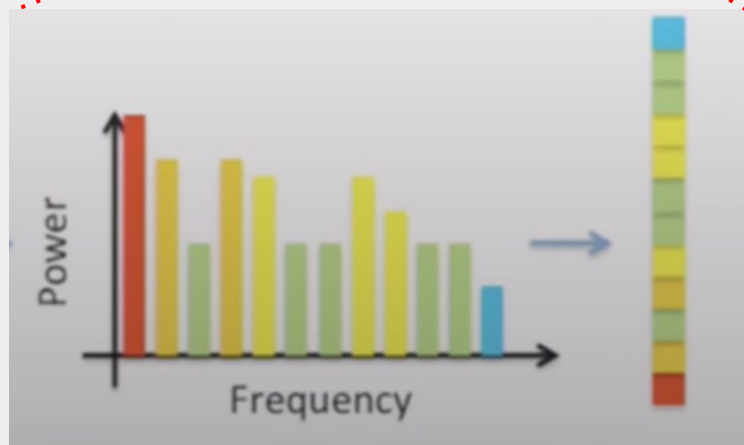


- Da li spektar govornog signala ostaje isti za svaki prozor?
- Izdelimo signal u vremenskom obliku na male prozore (npr. 20ms) i primenimo Furijeovu transformaciju na svaki.
- Snagu pojedinačnih frekvencija u tom prozoru smeštamo na spektrogram.

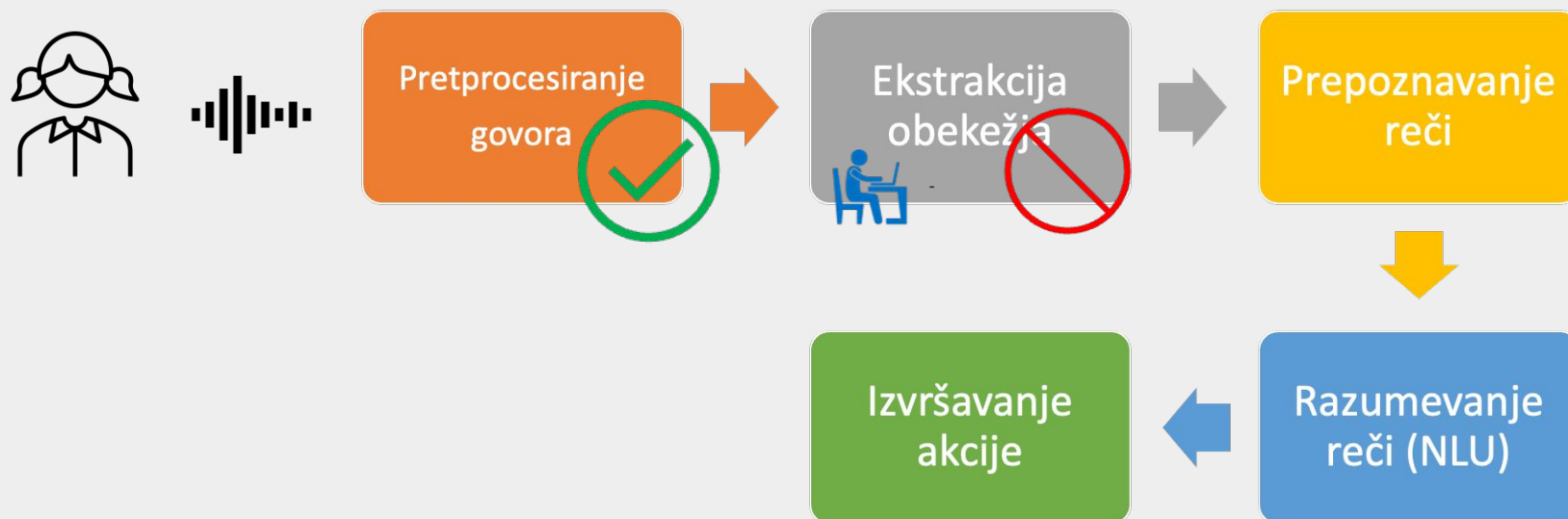
# Spektrogram



20ms



# Kako rade govorni asistenti?

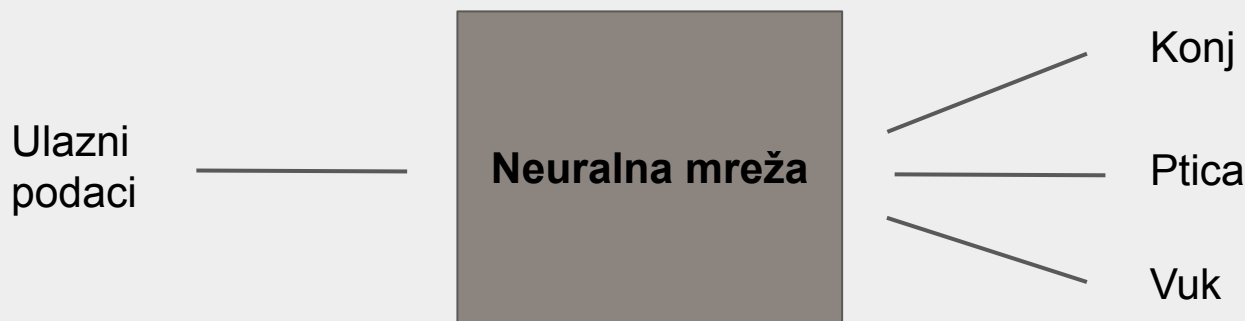


1. Uvod u obradu govora

2. Metode dubokog učenja

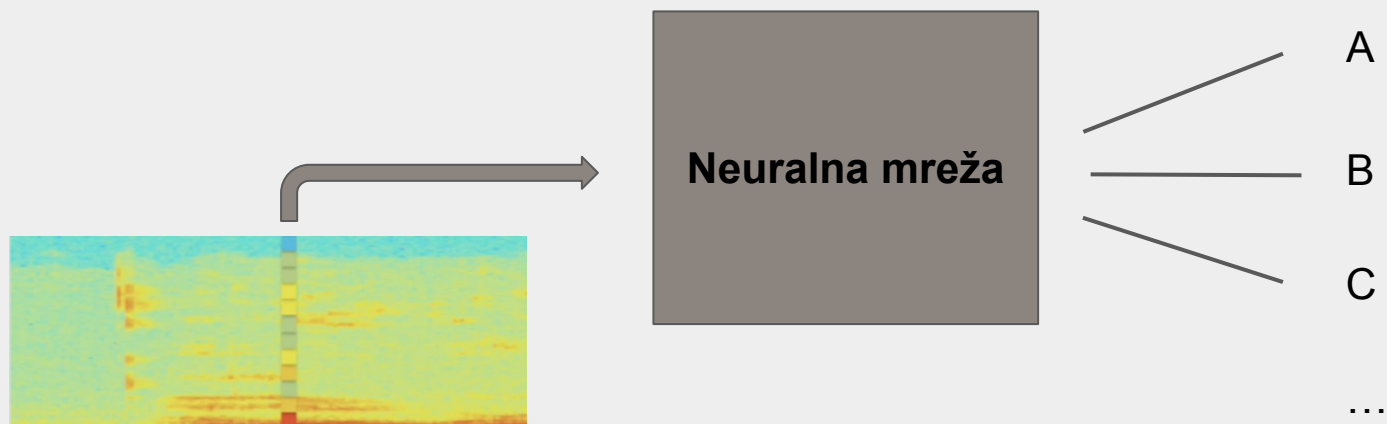
3. Prepoznavanje emocija

- **Duboko učenje** je grana mašinskog učenja koja ima za cilj da modeluje visoke nivoe apstrakcije u podacima koristeći se velikim brojem procesirajućih slojeva, sa ili bez kompleksih struktura i nelinearnim transformacijama.

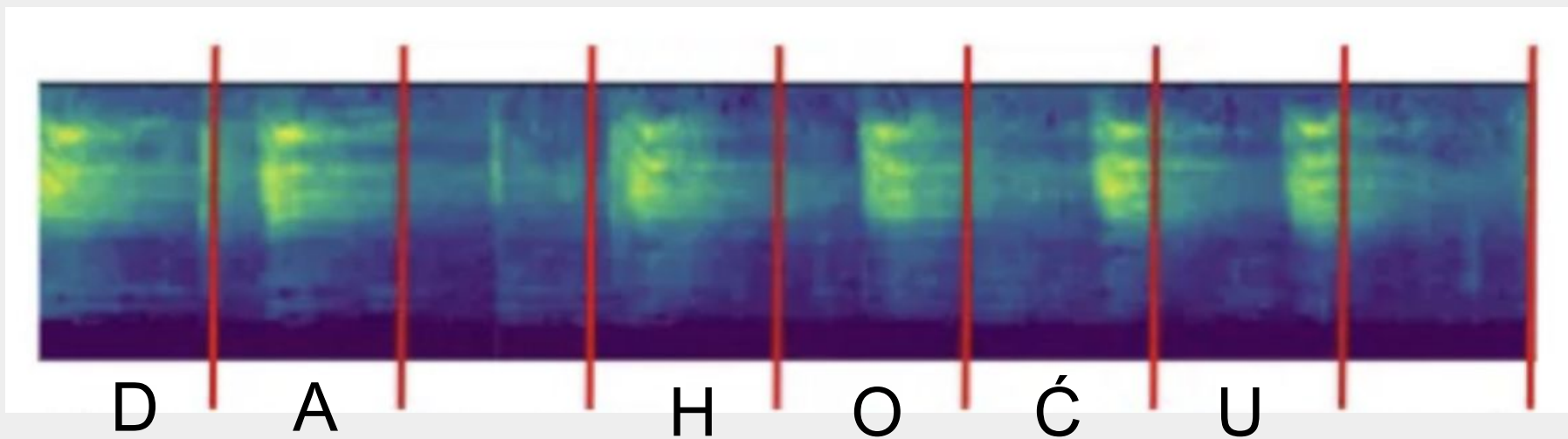


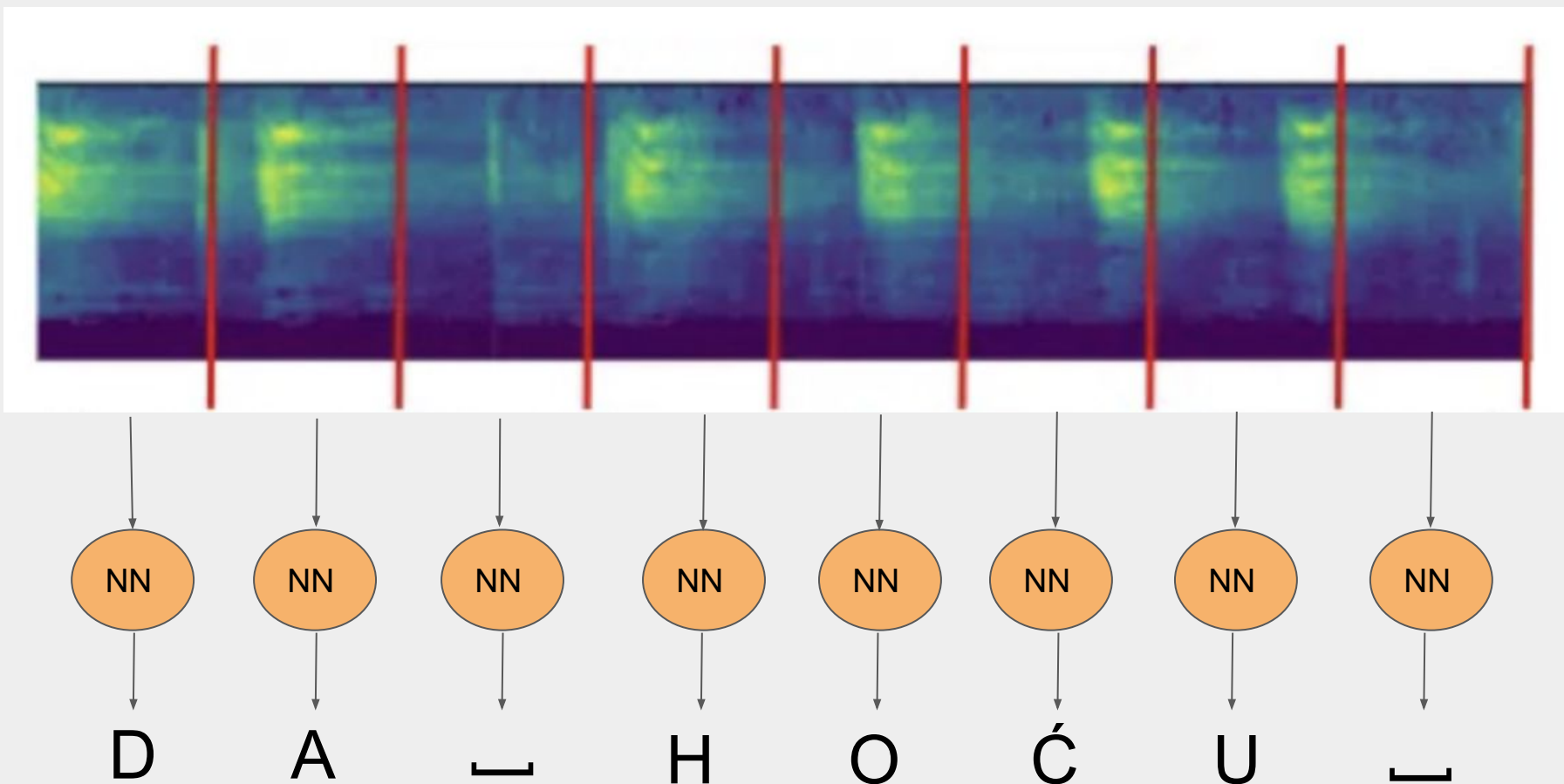
- Prepoznavanje govora, šta su
  - Ulazni podaci?
  - Klase?

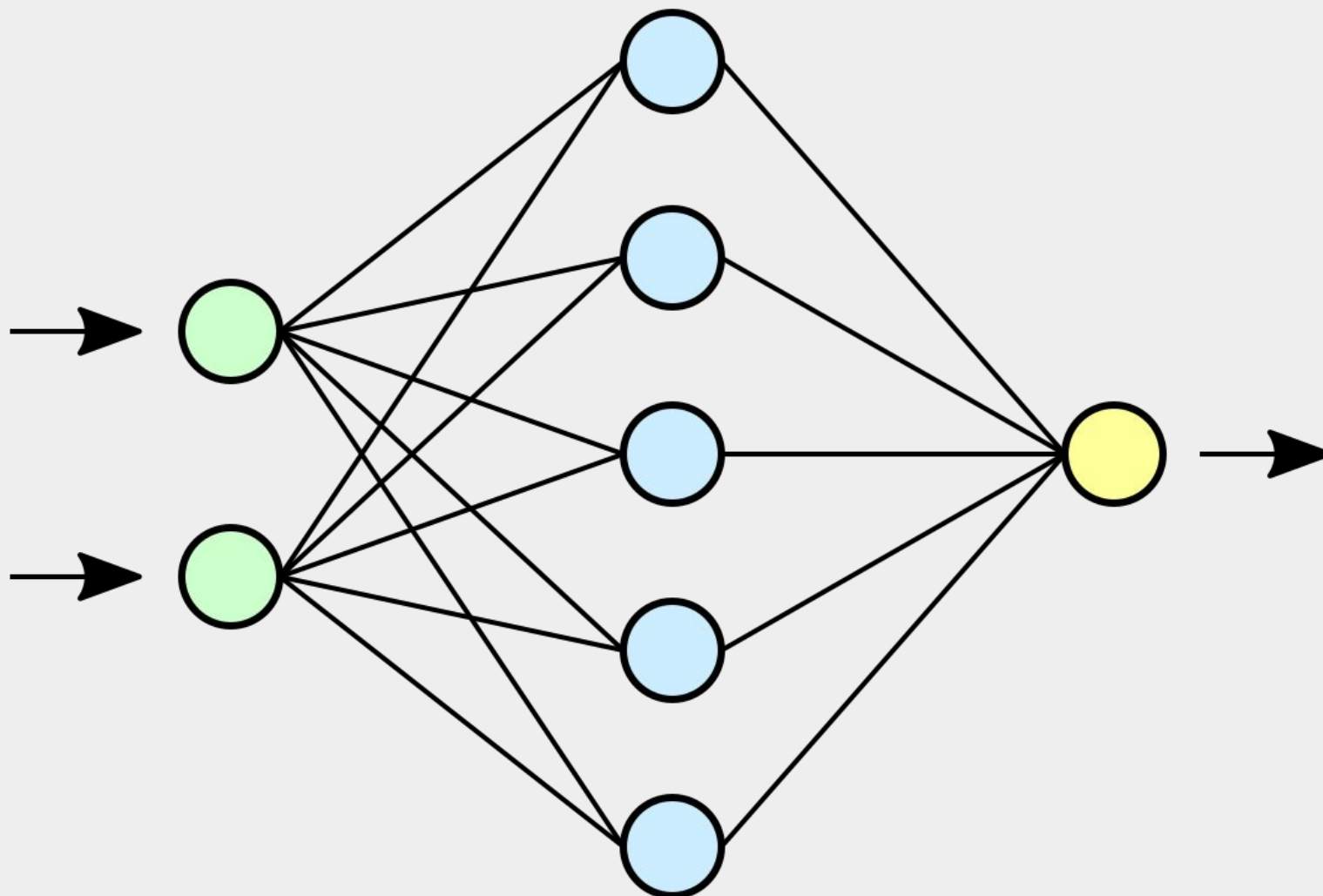
- Prepoznavemo slovo za svaki prozor spektrograma

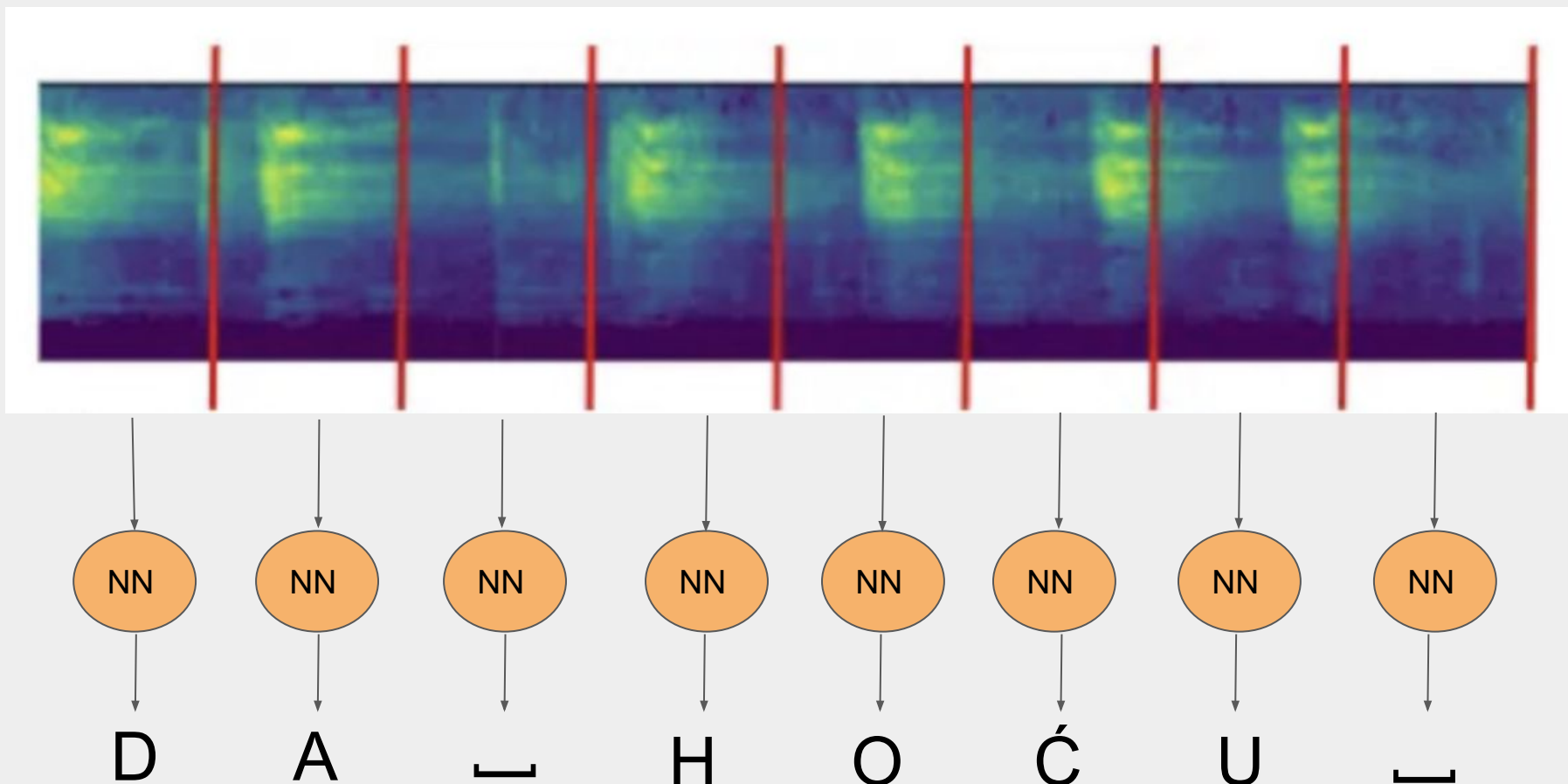






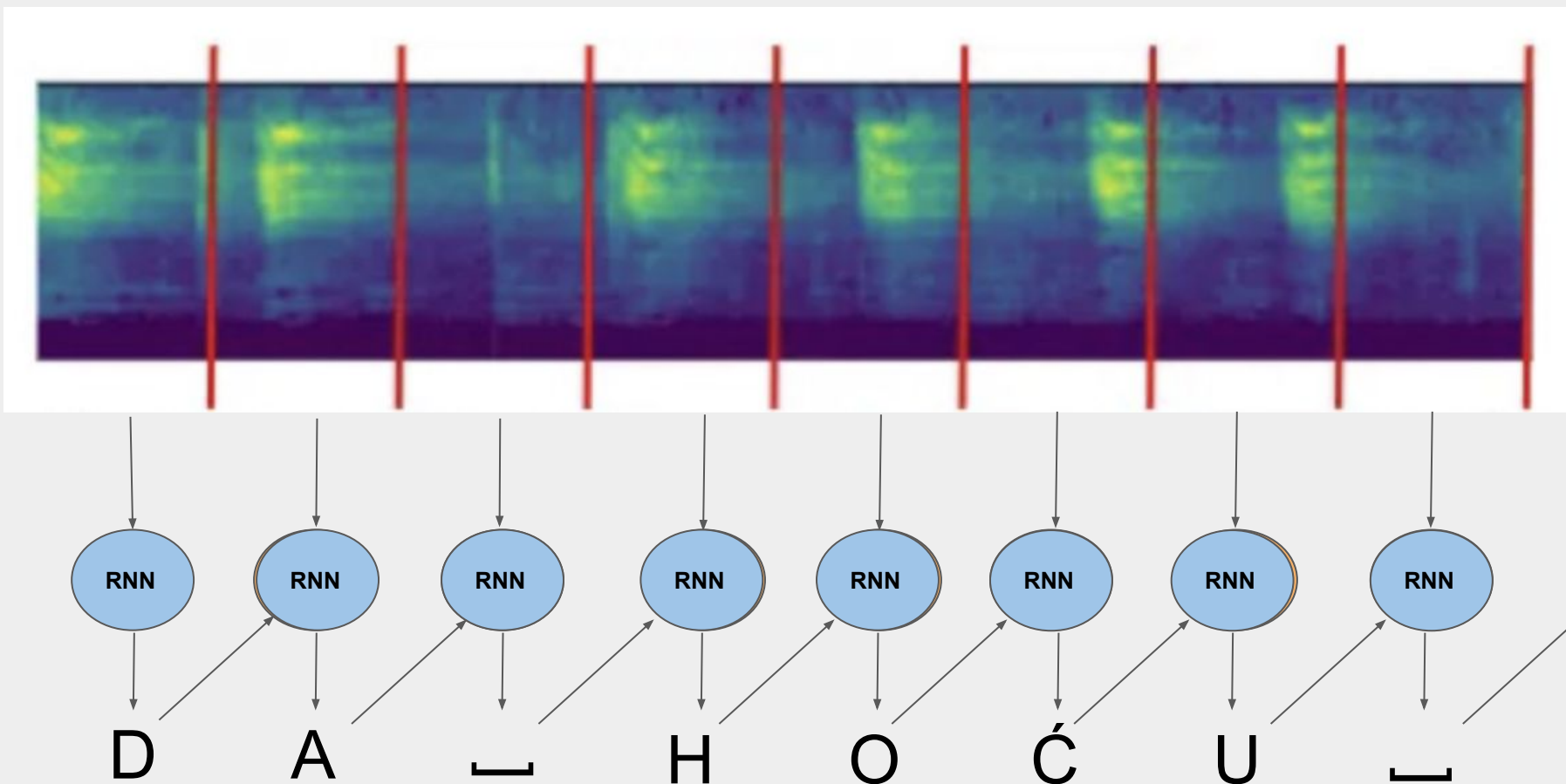


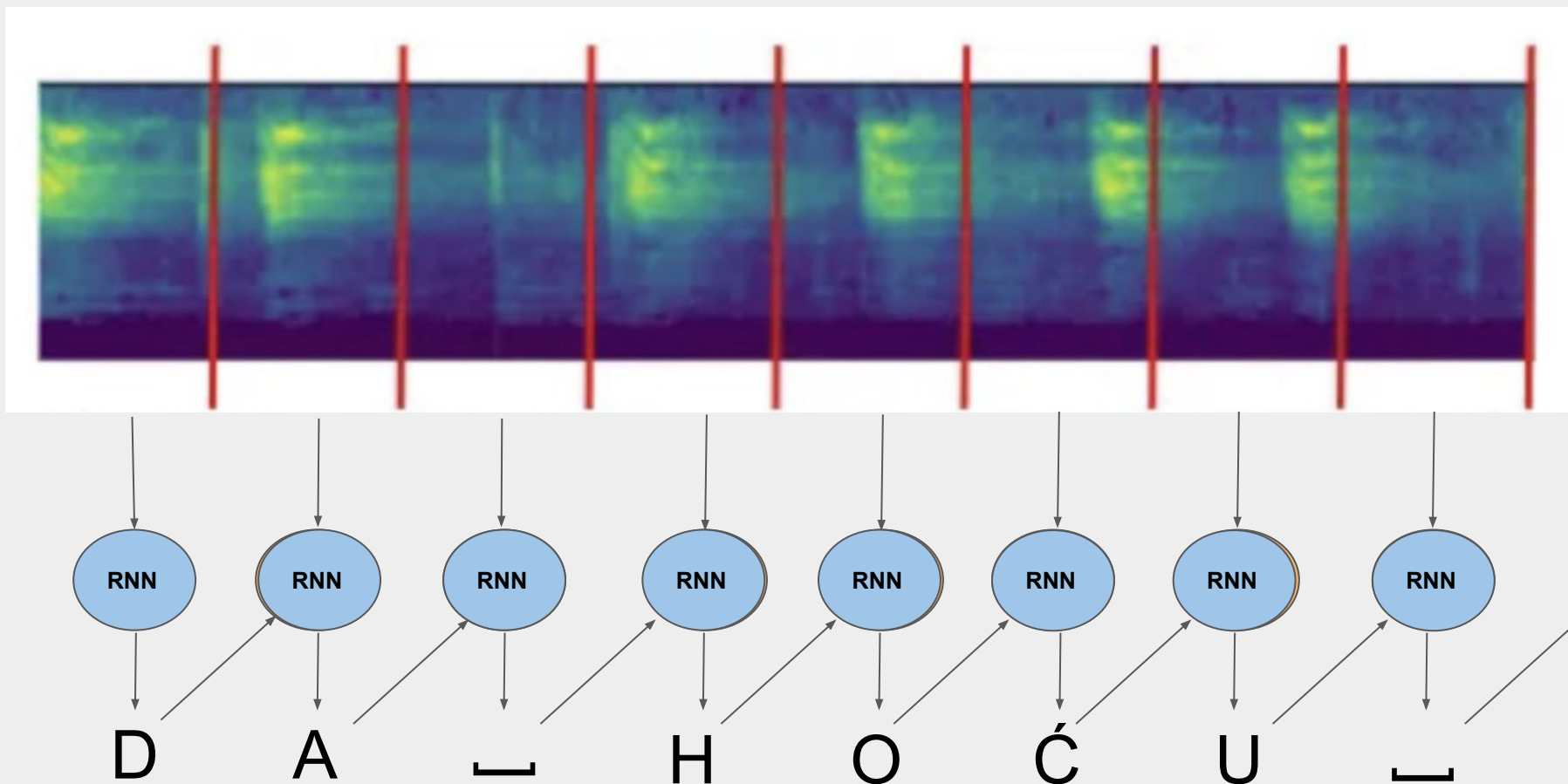




Da li koristimo sve dostupne informacije?

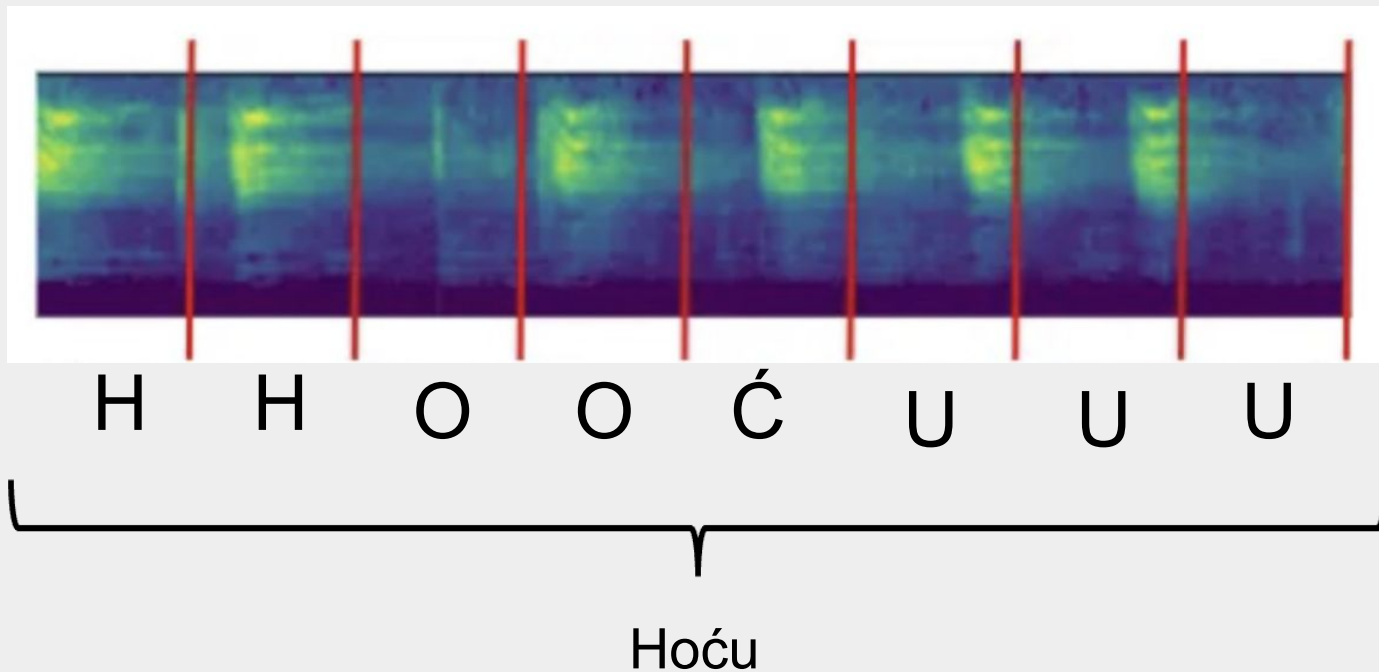
# Rekurentna neuralna mreža





Svaki prozor novo slovo? Hmm...

- Više uzastopnih vremenskih prozora odgovaraju istom slovu:

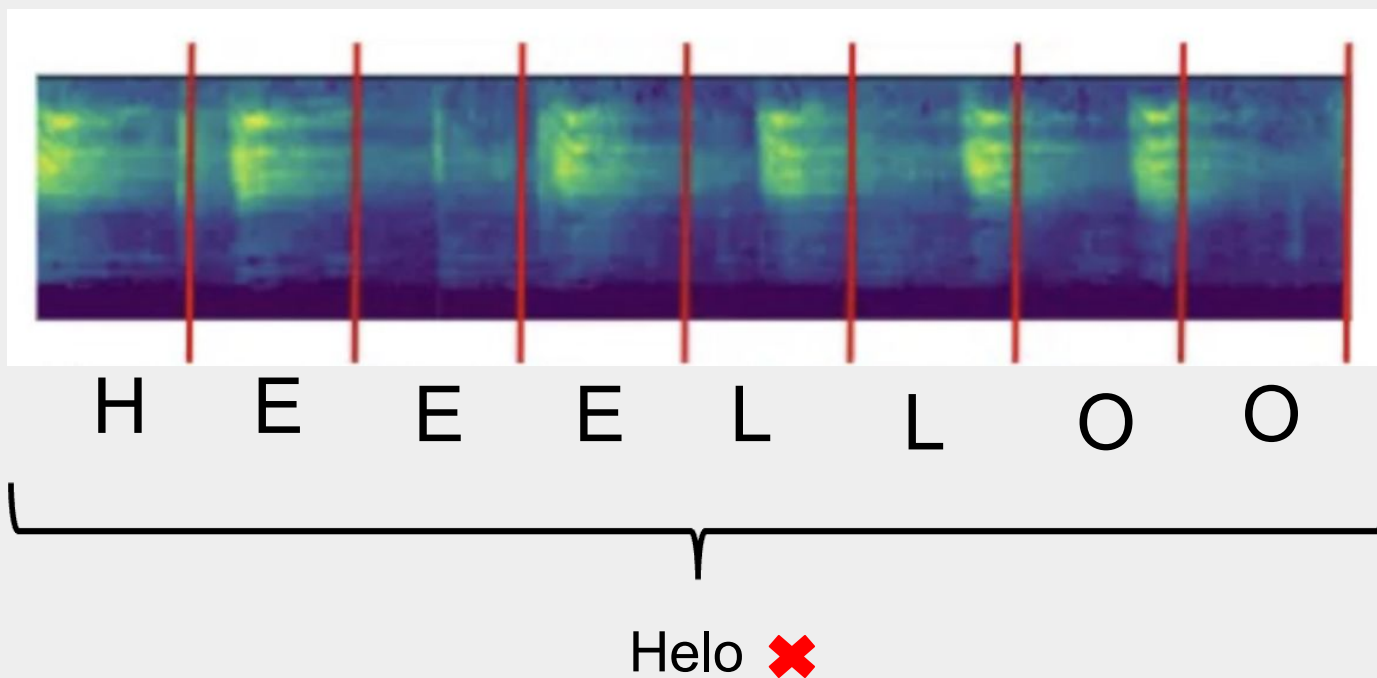


- Labeliranje pojedinačnih prozora je skupo, neefikasno i sklono greškama!
- Možemo li samo da sažmemo slova?

# Alignment problem (poravnanje)



- Dva ista slova za redom:





- Uvodimo *blank* karakter. Ovo nije isto što i razmak!
- Prvo sažmemo, pa uklonimo *blank*.

Potencijalni izlazi iz modela:

H \_ E \_ \_ L \_ L L \_ O -> HELLO V1

H H \_ E \_ L \_ L \_ O O -> HELLO V2

H \_ E E \_ L L L \_ O -> HELO

H \_ A \_ L \_ L \_ O O O -> HALLO

- Uvodimo *blank* karakter. Ovo nije isto što i razmak!
- Prvo sažmemo, pa uklonimo *blank*.

Potencijalni izlazi iz modela:

H \_ E \_ \_ L \_ L L \_ O -> HELLO V1

H H \_ E \_ L \_ L \_ O O -> HELLO V2

H \_ E E \_ L L L \_ O -> HELO

H \_ A \_ L \_ L \_ O O O -> HALLO

Kako ćemo da definišemo funkciju greške?



Potencijalni izlazi iz modela:

H \_ E \_ \_ L \_ L L \_ O -> HELLO V1

H H \_ E \_ L \_ L \_ O O -> HELLO V2

H \_ E E \_ L L L \_ O -> HELO

H \_ A \_ L \_ L \_ O O O -> HALLO

---

**CTC** (Connectionist temporal classification) **funkcija greške**:  $\max(P_m(\text{"HELLO"}))$

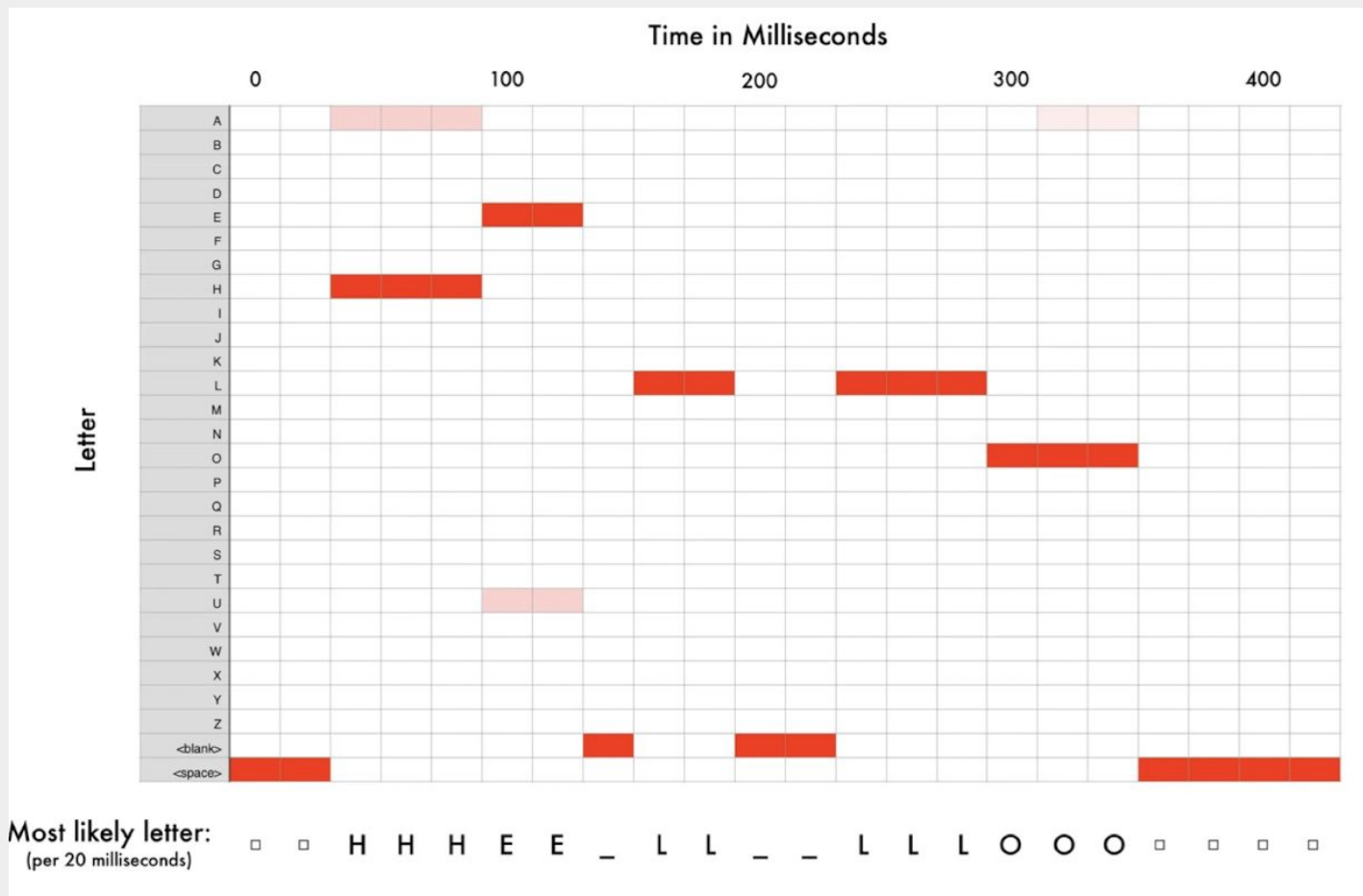
$$P_m(\text{"HELLO"}) = P(V1) + P(V2) + P(V3) + \dots$$

$$\text{gde je: } P(V1) = P(H) * P(\_) * P(E) * P(\_) * P(\_) \dots$$

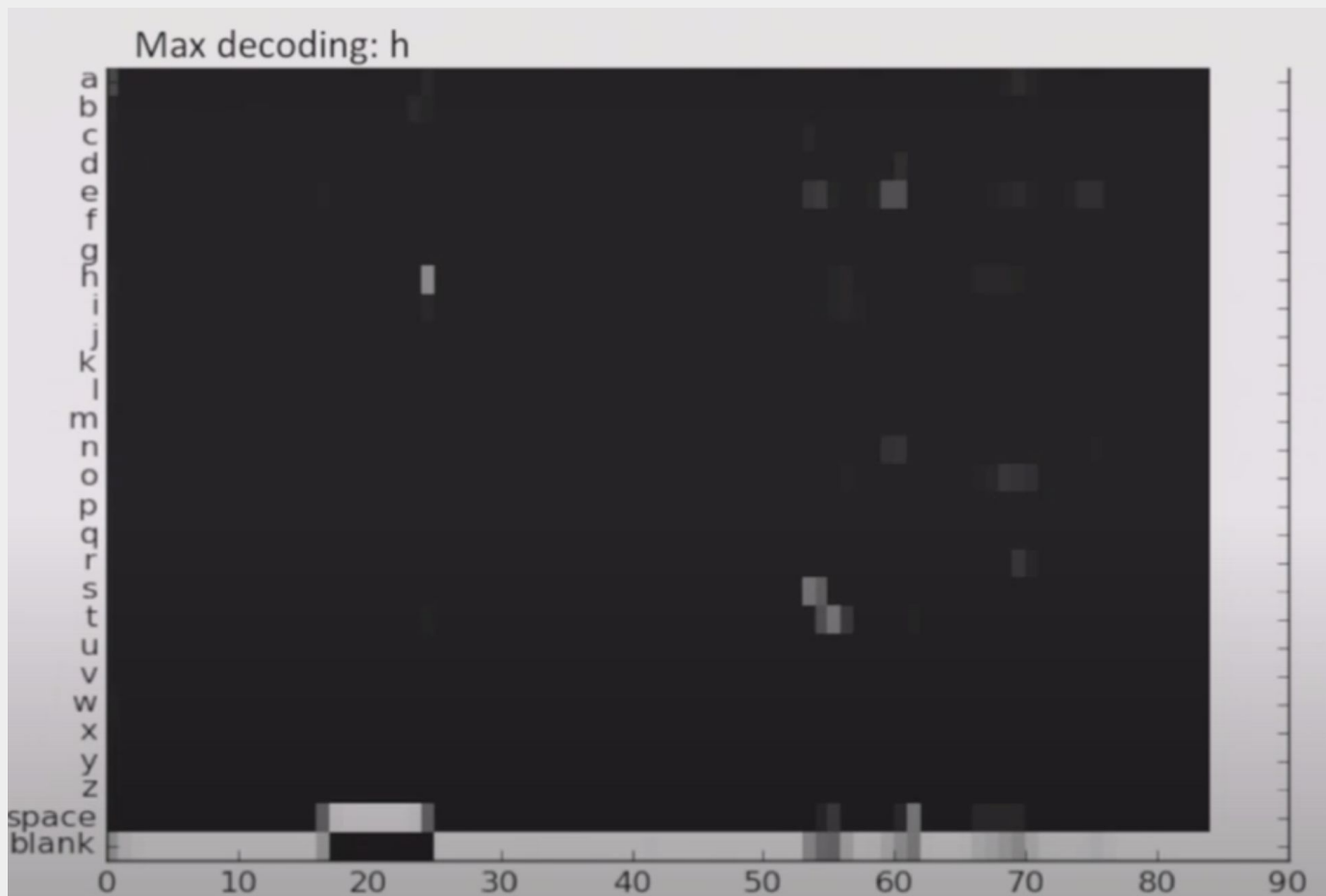
# Dekodovanje - Max Decoding



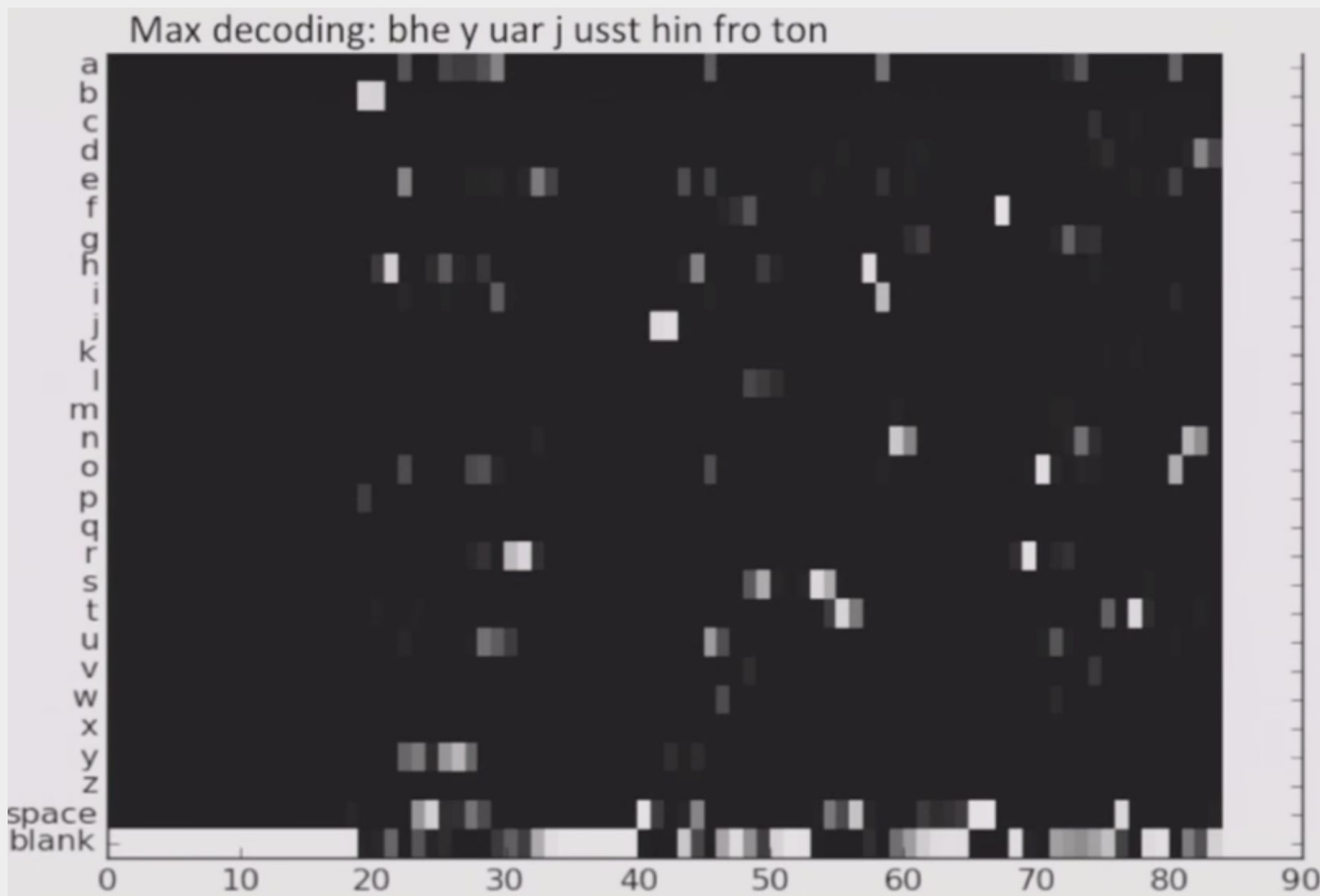
- Uvek biraj slovo sa najvećom verovatnoćom u tom vremenskom prozoru.



- Nakon 300 iteracija

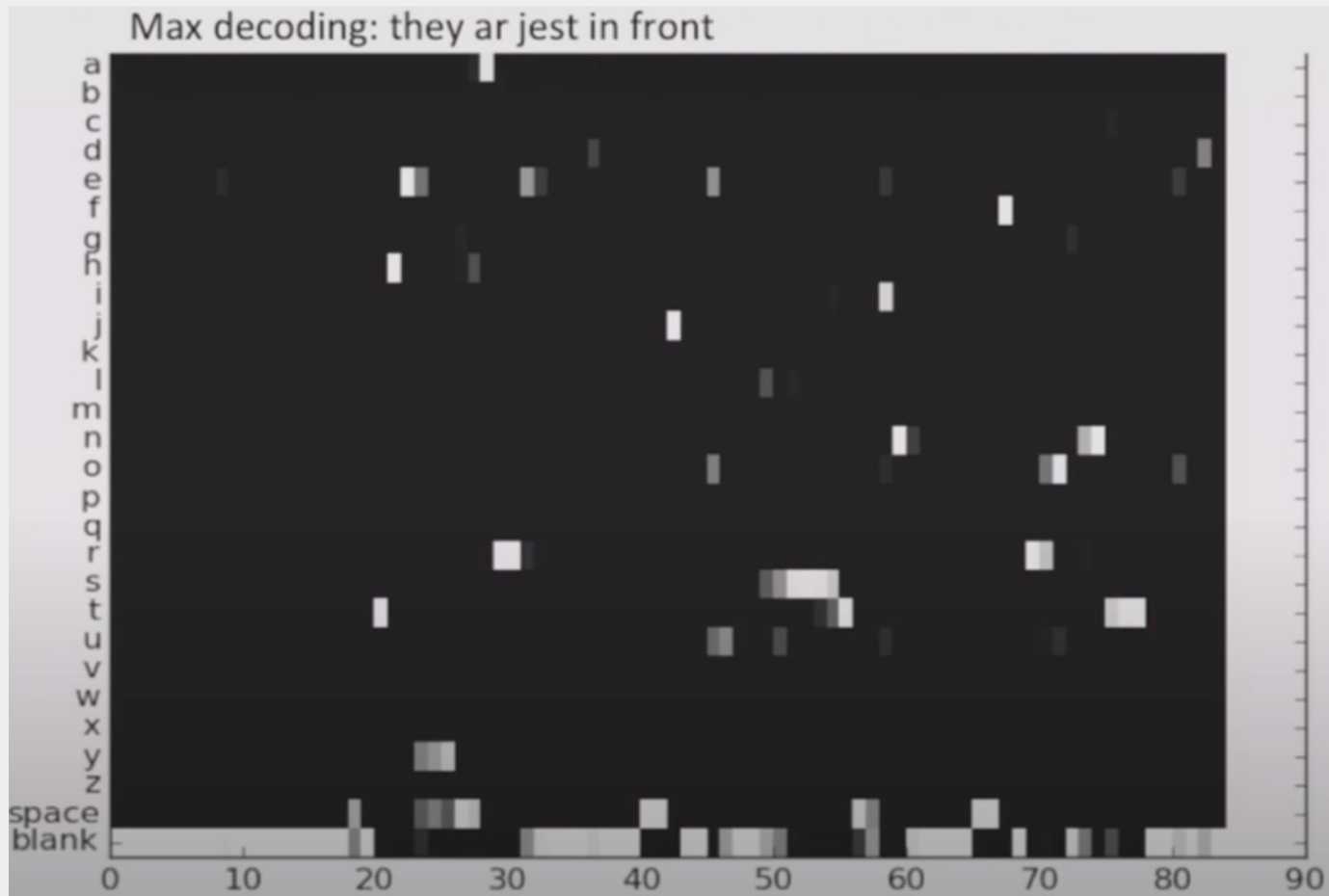


- Nakon 1500 iteracija



- Nakon 5500 iteracija

Izgovoreno: “There just in front”



- Dobre sekvence su često odbačene prerano. Sistemi za prepoznavanje govora nemaju dovoljno bogate podatke za trening. (Vlastite imenice, retke reči...)
- Možemo da iskoristimo moć jezičkih modela (NLU), trenirani su na više podatak i mogu bolje da procene verovatnoću retkih sekvenci.

Koristimo jezički model  $P_{LM}$ :

$$P(C_t) = P_m(C_t) * P_{LM}(C_t | C_{1:t-1})$$

Do sada: **LO**

$$\text{Next: O? } P(\text{LO}\mathbf{V}) = P_m(\mathbf{V}) * P_{LM}(\mathbf{V} | \text{LO})$$

$$\text{Next: M? } P(\text{LO}\mathbf{M}) = P_m(\mathbf{M}) * P_{LM}(\mathbf{M} | \text{LO})$$

Do sada: **LA**

$$\text{Next: O? } P(\text{LA}\mathbf{V}) = P_m(\mathbf{V}) * P_{LM}(\mathbf{V} | \text{LA})$$

$$\text{Next: M? } P(\text{LA}\mathbf{M}) = P_m(\mathbf{M}) * P_{LM}(\mathbf{M} | \text{LA})$$

Previše kombinacija -> Beam search

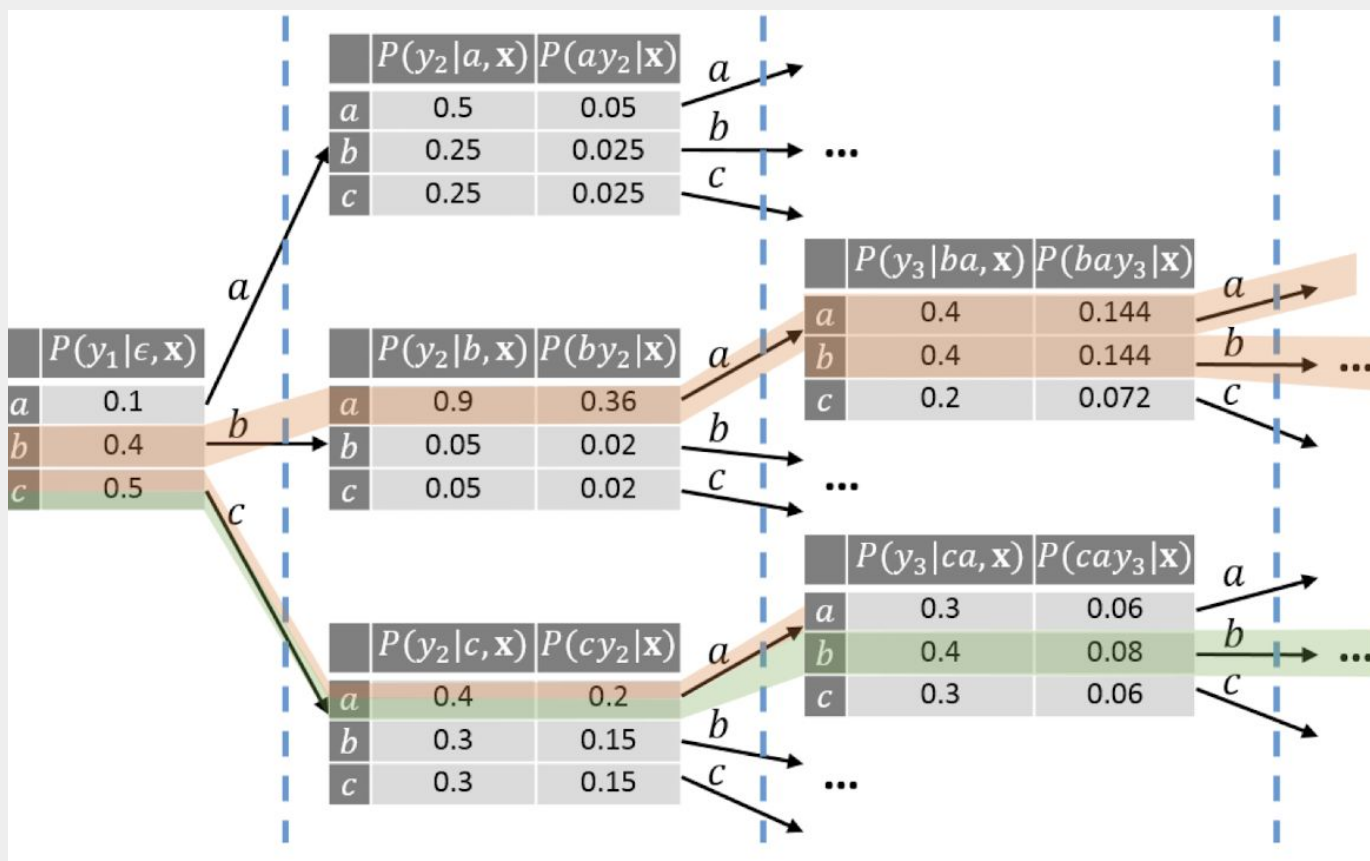


# Dekodovanje - Pretragom po snopu (Beam search)



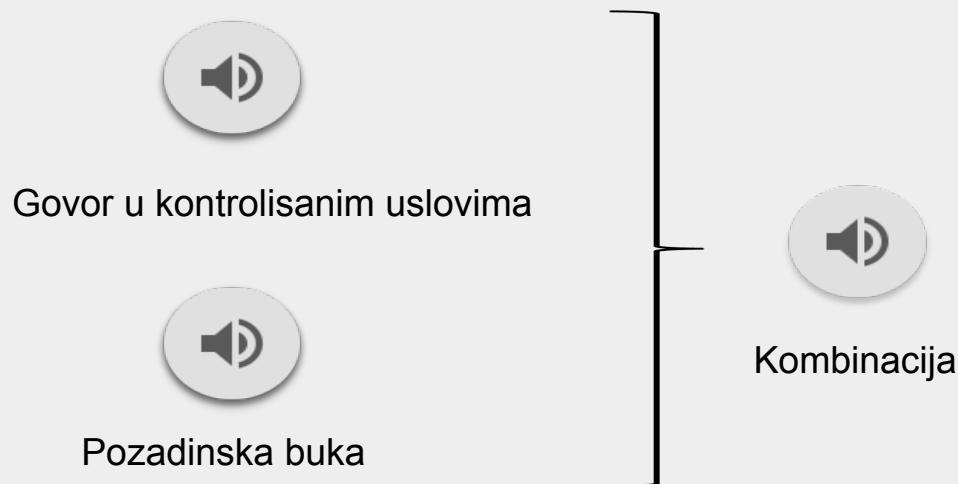
Koristimo jezički model  $P_{LM}$ :

$$P(C_t) = P_m(C_t) * P_{LM}(C_t | C_{1:t-1})$$



- Najčešće se koristi „čitani govor” iz izvora kao što su *Street Journal* (novine) i *LibriSpeech* (audio knjige).
- **Kvalitet podatak** je **izuzetno važan** kod metoda dubokog učenja, neadekvatni podaci dovešće do niske performanse modela u realnoj upotrebi.
- Ovi primeri ne odgovoraju većini svakodnevnih situacija ili specifičnih primena.
- Lombard efekat.
- Akcenti, pol, zamuckivanja, godište, šumovi...

- Veštačko formiranje skupa podataka veće raznovrsnost.
- Povećanje/smanjenje brzine, dodavanje šuma...
- Augmentacijom podataka možemo povećati skup podataka nekoliko puta.
- Nekad je jednostavnije uticati na podatke nego na robustnost modela.



- Često se umesto prepoznavanja slova, prepoznaju **fonemi**.
- Fonemi su eksperimentalno određene zvučne jedinice jezika.

Dekompozicija na:

*foneme*

*slova*

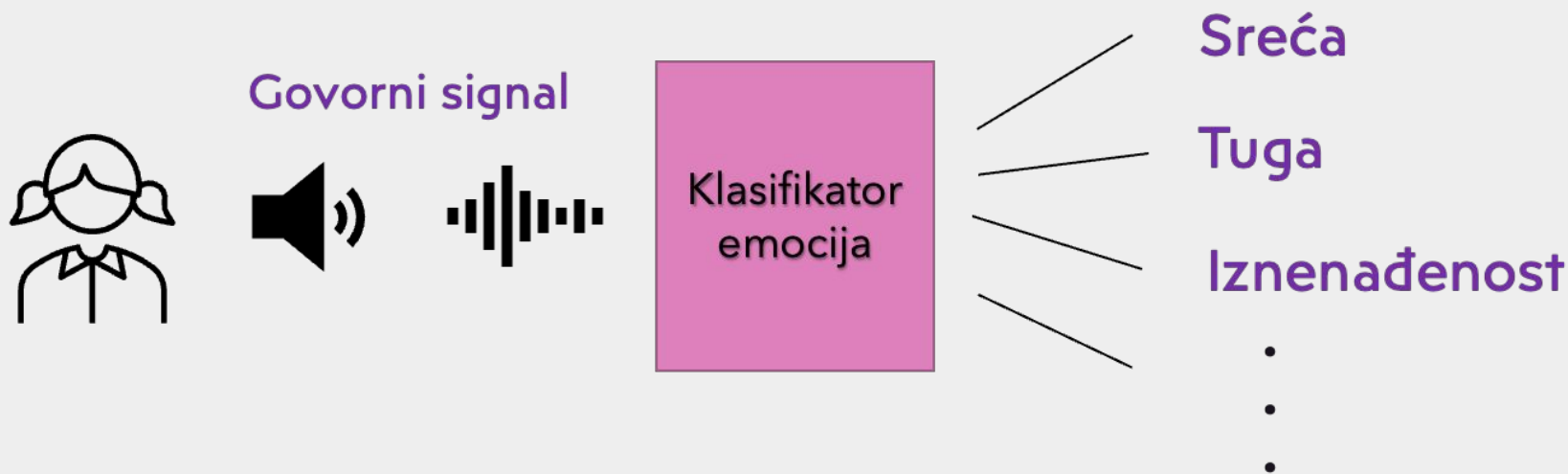
Hello -> ['HH' 'AH' 'L' 'OW']

Hello -> ['H' 'E' 'L' 'L' 'O']

	Phone Label	Example		Phone Label	Example		Phone Label	Example
1	iy	beet	22	ch	choke	43	en	button
2	ih	bit	23	b	bee	44	eng	Washington
3	eh	bet	24	d	day	45	l	lay
4	ey	bait	25	g	gay	46	r	ray
5	ae	bat	26	p	pea	47	w	way
6	aa	bob	27	t	tea	48	y	yacht
7	aw	bout	28	k	key	49	hh	hay

1. Uvod u obradu govora
2. Metode dubokog učenja
3. Prepoznavanje emocija

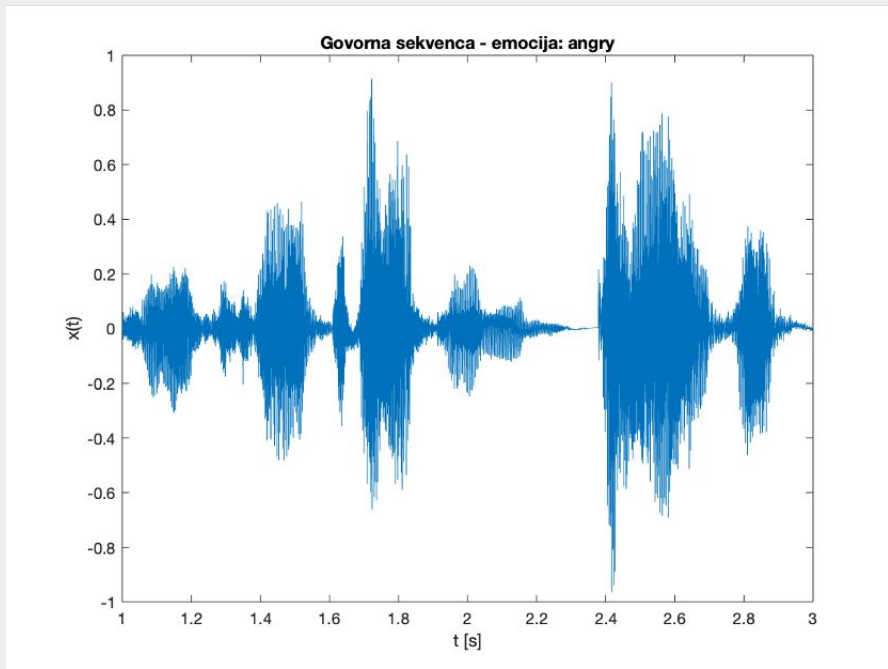
- Klasifikacioni problem.
- Prepoznavanje emocija može biti subjektivno.
- Primene
  - Govorna podrška (analiza i unapređenje)
  - Automatski intervju asistenti
  - Razumevanje budućih akcija
  - Terapije i mentalno zdravlje



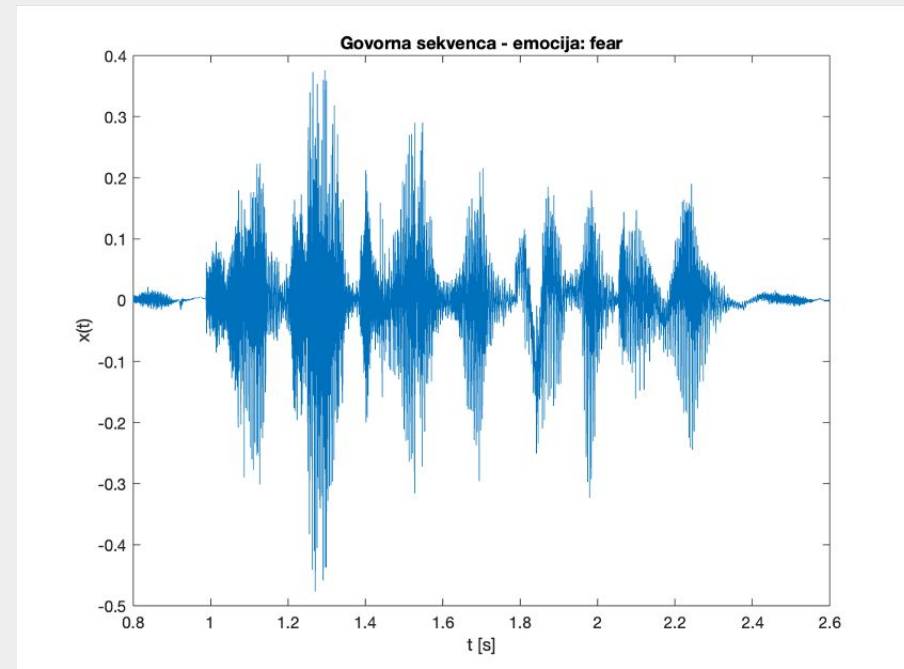
- RAVDESS baza podataka
- 24 govornika iz istog govornog područja - glumci
- 2 rečenice
- 8 emocija

Sad    Angry    Disgust    Fear    Surprise    Calm    Neutral    Happy

- Odsecanje tišine postavljanjem praga snage

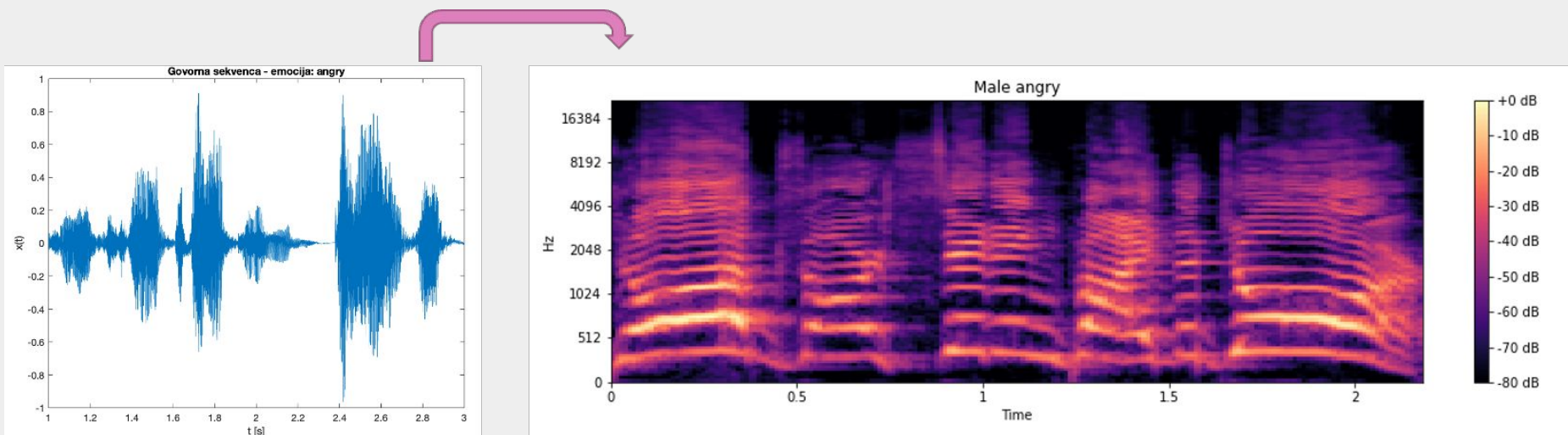


Ljutnja



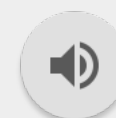
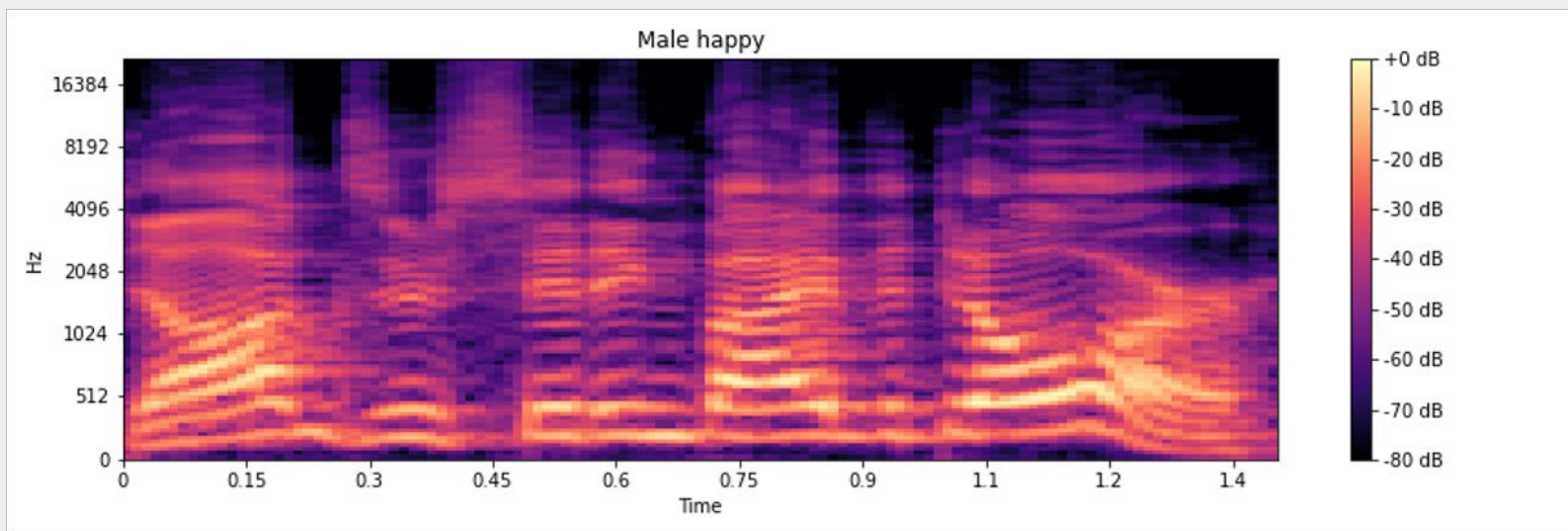
Strah



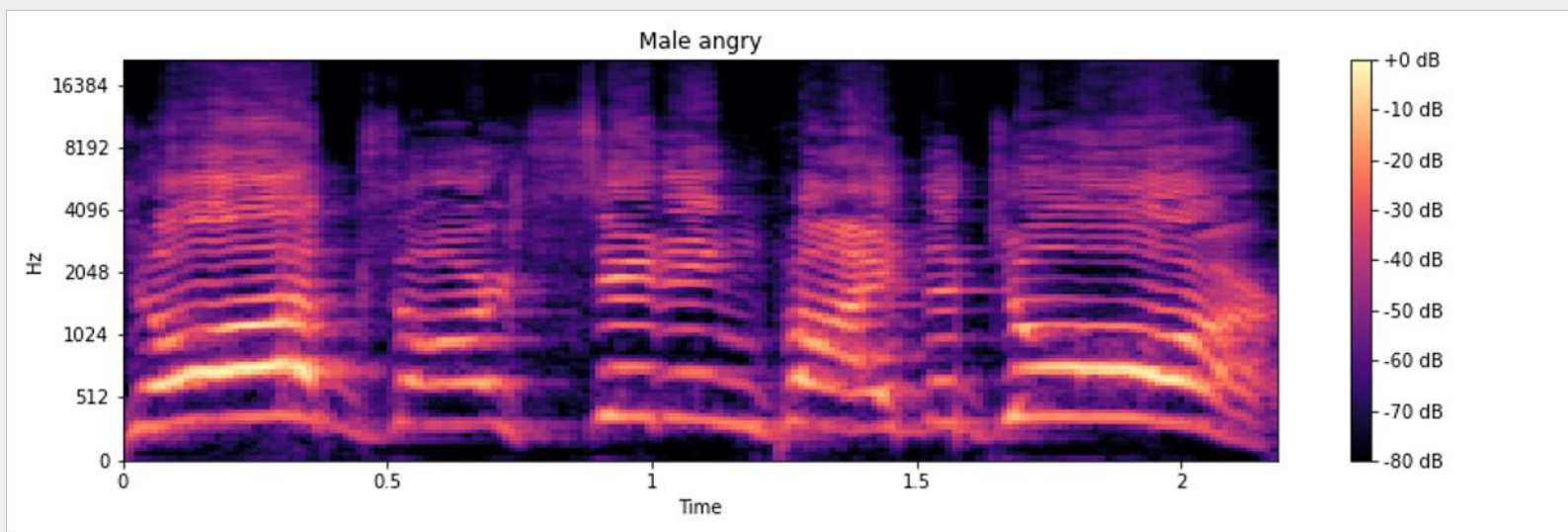


- Mel skala - ljudsko uho teže razlikuje tonove na višim frekvencijama.
- Eksperimentalna frekvencijska skala (mel skala)

# Spektrogram emocija

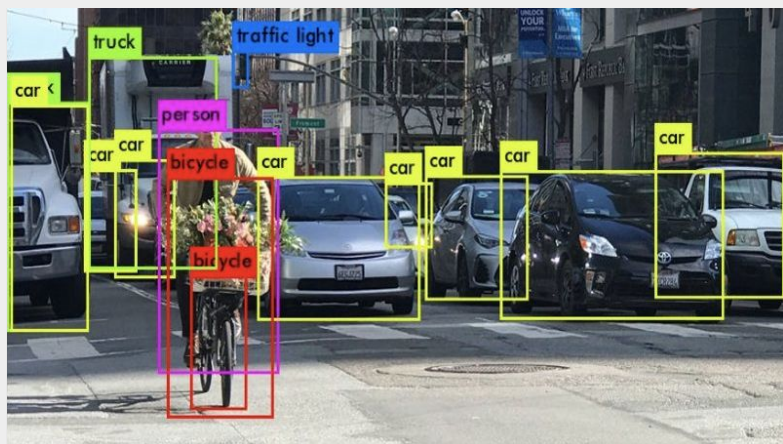


Sreća

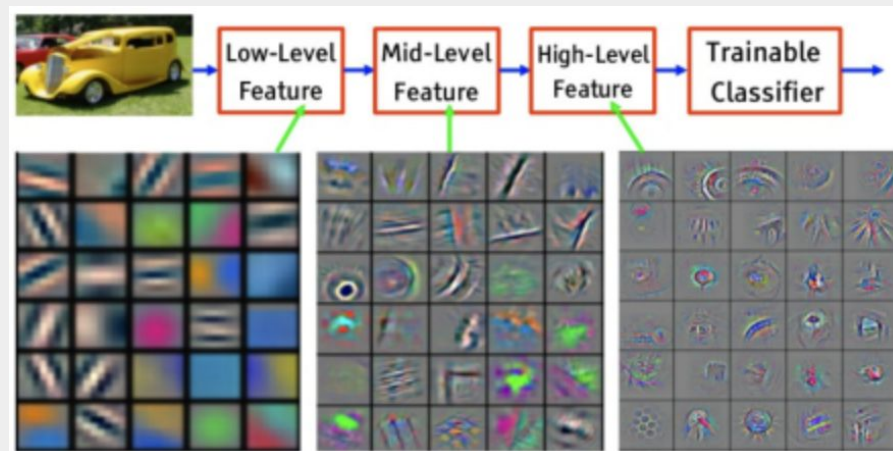


Ljutnja

- Specijalni tip neuralnih mreža specijalizovan za izvlačenje podataka iz mrežastih struktura koa što su slike

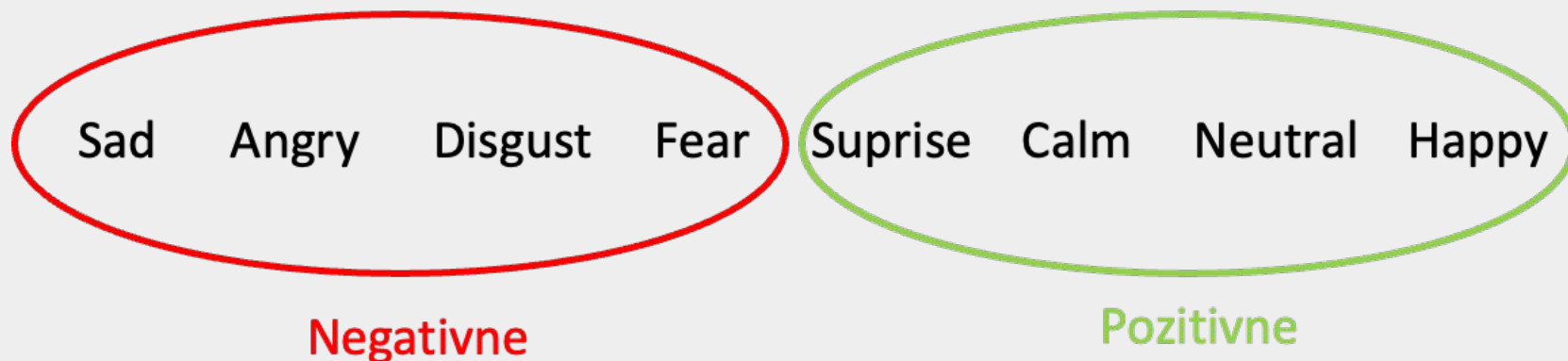


CNN za detekciju objekata



Specijalizacija za konkretna obeležja kroz slojeve

- Binarna klasifikacija



- Tačnost: **84%**

- Po tipu - 8 klasa

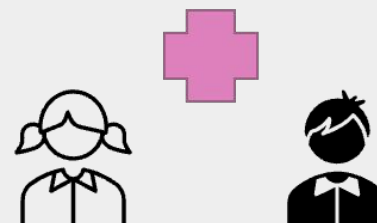
Sad    Angry    Disgust    Fear    Suprise    Calm    Neutral    Happy

- Tačnost: **52%**, Top 3: **88%**

- Po tipu i polu - 16 klasa

Sad    Angry    Disgust    Fear    Suprise    Calm    Neutral    Happy

- Tačnost: **50%**, Top 3: **78%**





Sad   Angry   Disgust   Fear   Surprise   Calm   Neutral   Happy



Labela

Model

Iznenadjena

Srećna

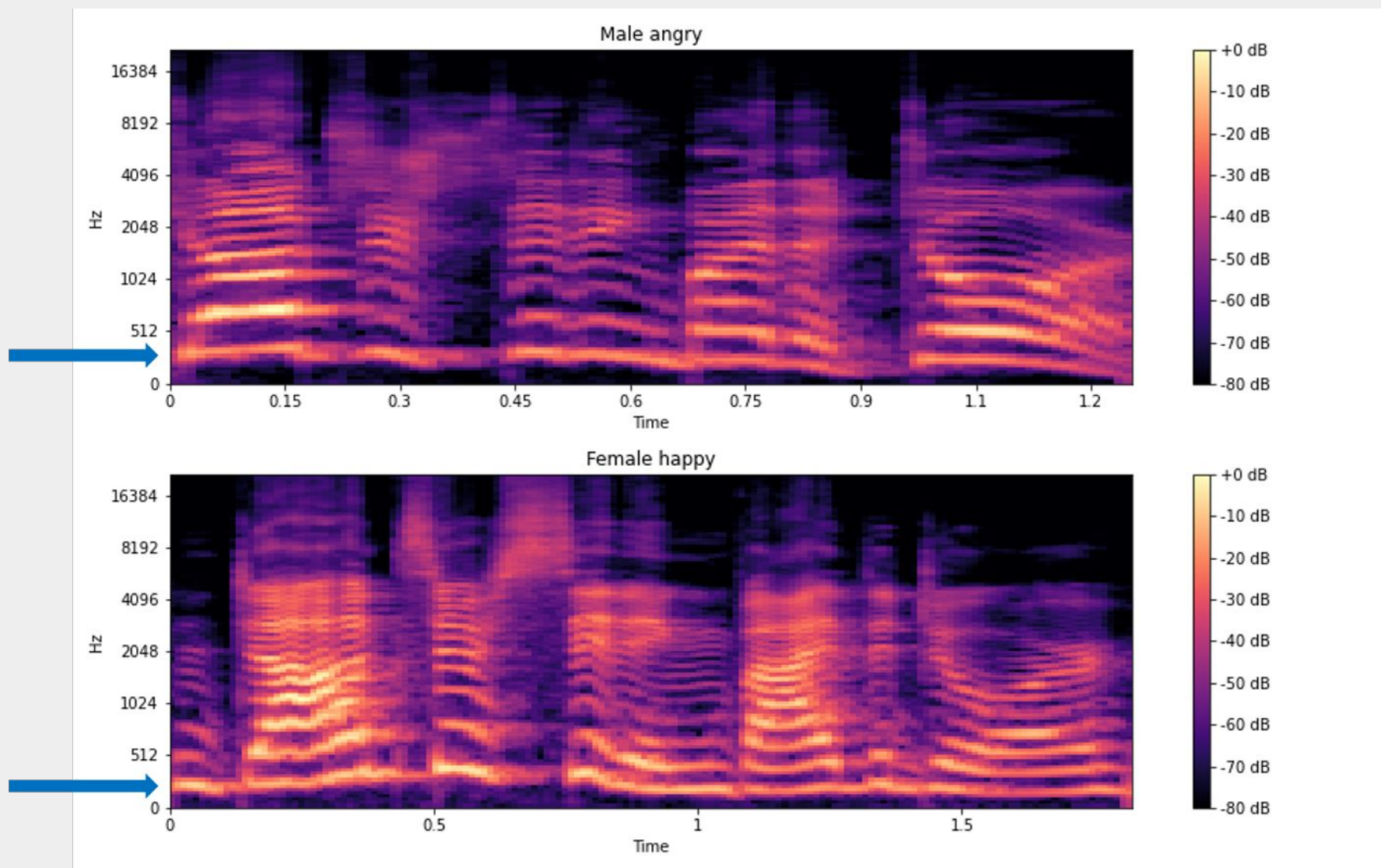
Zgrožen / Ljut

Zgrožen / Uplašen ...

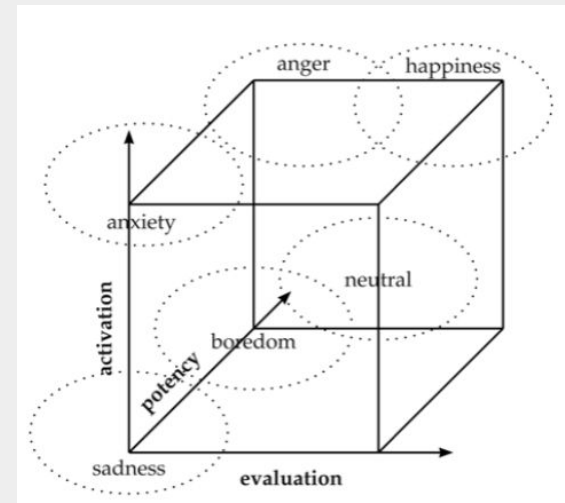
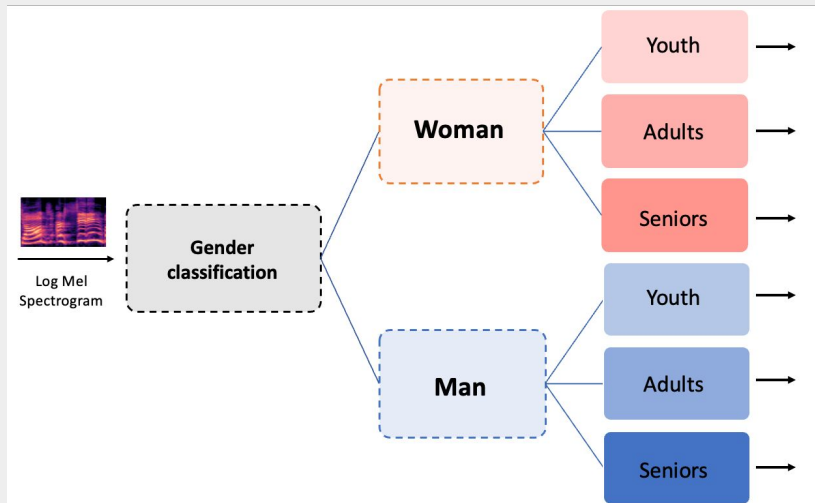
		Actor intended emotion							
		Neutral	Calm	Happy	Sad	Angry	Fearful	Disgust	Surprise
Rater chosen emotion	Neutral/Calm	86.6	69.9	14.25	17.12	4.03	4.5	4.36	7.03
	Happy	0.63	17.27	68.44	1.48	0.23	0.59	0.59	6.56
	Sad	4.65	6.06	2.29	60.85	1.02	6.58	8.65	0.76
	Angry	3.82	1.02	1.79	2.9	81.32	4.79	6.48	2.78
	Fearful	0.63	0.66	1.67	9.64	1.39	70.71	2.31	2.22
	Disgust	1.15	1.46	0.78	3.09	8.37	1.81	69.77	3.28
	Surprise	0.28	0.33	7.88	0.69	1.2	7.76	4.13	72.29
	None	2.26	3.3	2.9	4.24	2.45	3.26	3.72	5.07



# Ljut muškarac = Srećna žena?



- Velika neponovljivost i nekonzistentnost labela među ispitanicima.
- Dodavanje informacija o izgovorenom tekstu ili facijalnim ekspresijama.
- Pojedinačni modeli za različite polove i starosna doba.
- Pretstavljane emocija u prostoru aktivacija-valentnost-dominantnost.





Hvala na pažnji!

Pitanja?

$$f_N(t) = \sum_{n=-N}^N c_n e^{in\omega t}$$

$$m = 2595 \log_{10} \left( 1 + \frac{f}{700} \right)$$



- Zdravo, ovo je neki tekst
  - On može imati i uvučene teze
- ... a i redovne teze
- ... šta god ti duša ište