

Језички модели (за српски језик)

Обучавање (различитих) модела

**Композитни модели
(за евалуацију и генерисање
текста)**



др Михаило Шкорић

Наставница пита
да ли је реченица
исправна.

Перице, да ли је ова
реченица граматички
исправна?



ЧАС СРПСКОГ
-садашњост-

Од деце се очекује
буловски одговор
(јесте/није
true/false)

Само мало да
размислим...

јесте

није, није



ФОРМАЛНА ГРАМАТИКА

Господине, па
ви сте робот!

Ово је
реченица
мог језика.

true

На исти начин функционише формална граматика.

Формална граматика српског језика нам говори да ли је
неки улазни текст на српском, тј. да ли припада језику.

ФОРМАЛНА ГРАМАТИКА

Господине, па
ви сте робот!

Quoi???

false

Формална граматика неког другог језика ради то исто,
али за тај други језик (нпр. француски).

ФОРМАЛНА ГРАМАТИКА

Па ви сте
робот,
monsieur!

Euuuh...

false

Али шта се дешава ако је улазни текст мешовитог порекла?



ФОРМАЛНА ГРАМАТИКА

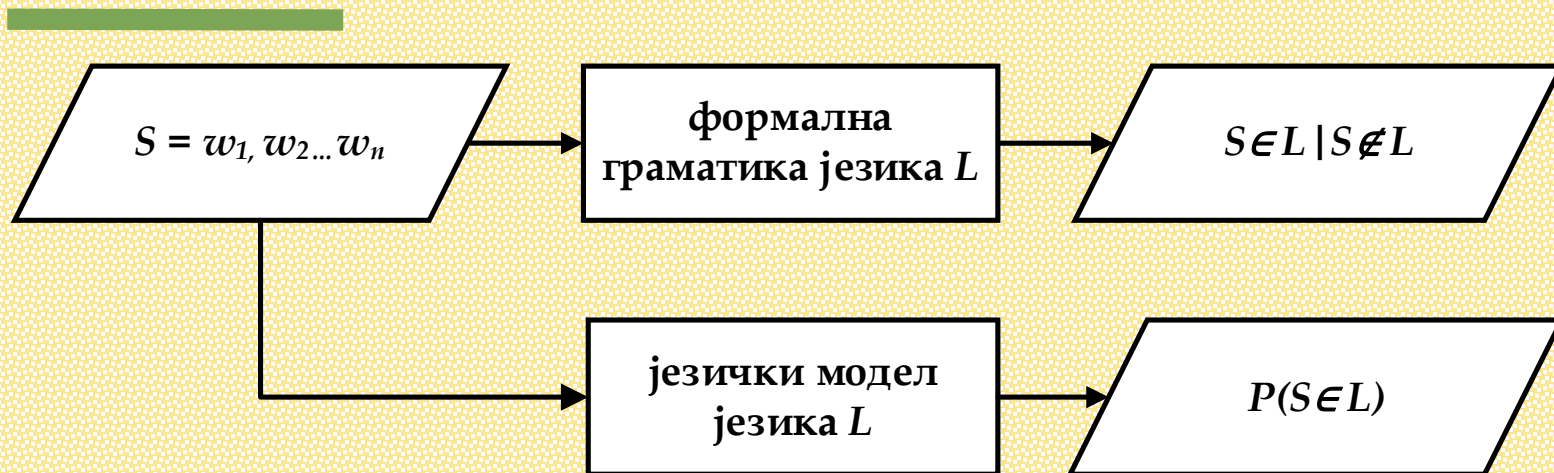
Па ви сте
робот,
monsieur!

ХМММ...

false

Услед ограничења буловског одговора,
све граматике ће, вероватно, рећи да је неисправан.

Језички модели? Граматике? Псеудограматике?



Као одговор на тај проблем, намеће језички модел.

Језички модел је расподела
вероватноће над нискама текста.

ЈЕЗИЧКИ МОДЕЛ

Па ви сте
робот,
monsieur!

Звучи
познато. 57%
шансе да је
ово српски.

0.57

Он одговор на питање припадности дају у виду вероватноће.

За необичајене ниске текста нуди ниску вероватноћу.

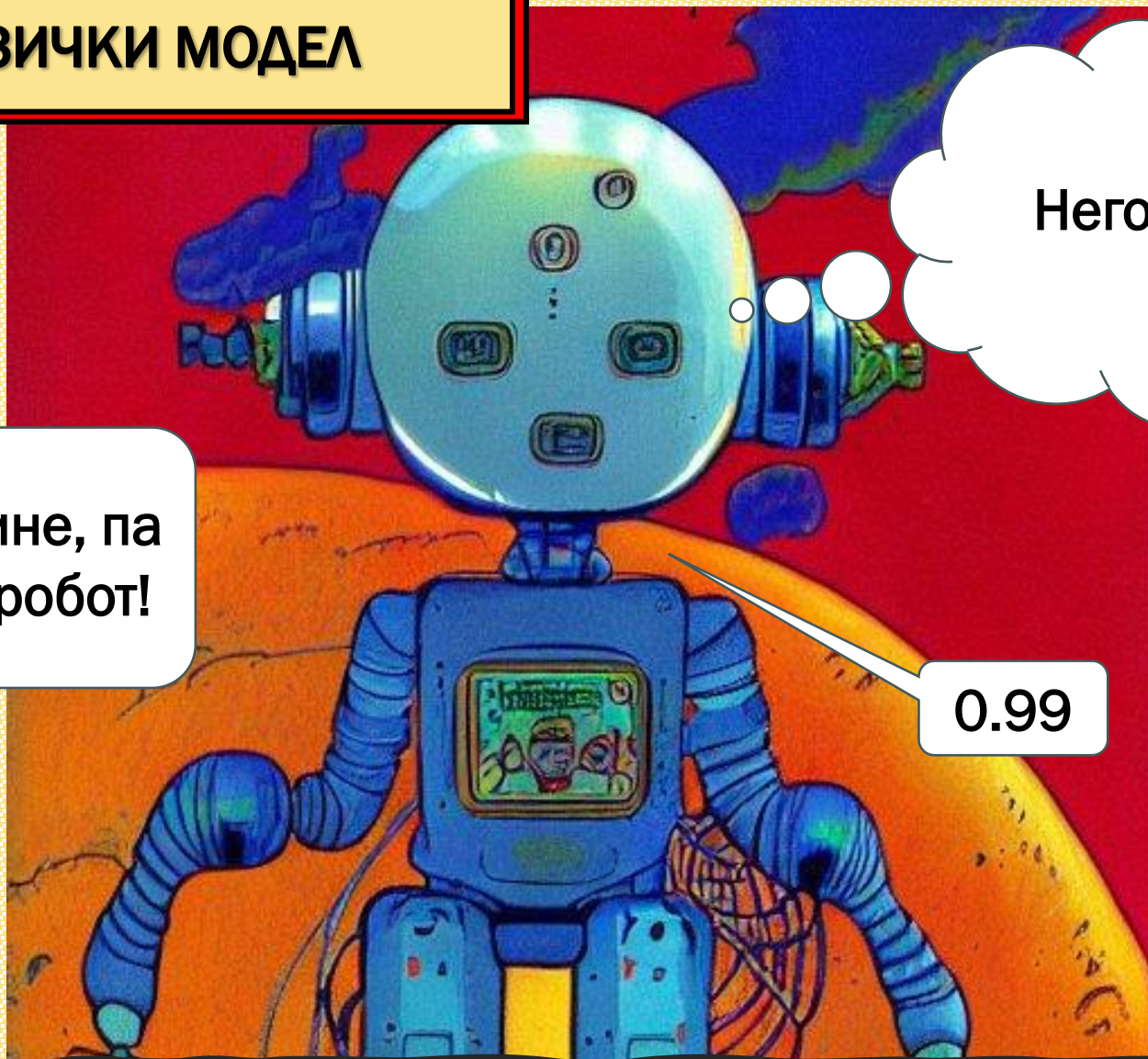
ЈЕЗИЧКИ МОДЕЛ


Господине, па
ви сте робот!

Него шта!

0.99

А за уобичајене, високу вероватноћу да је текст исправан.





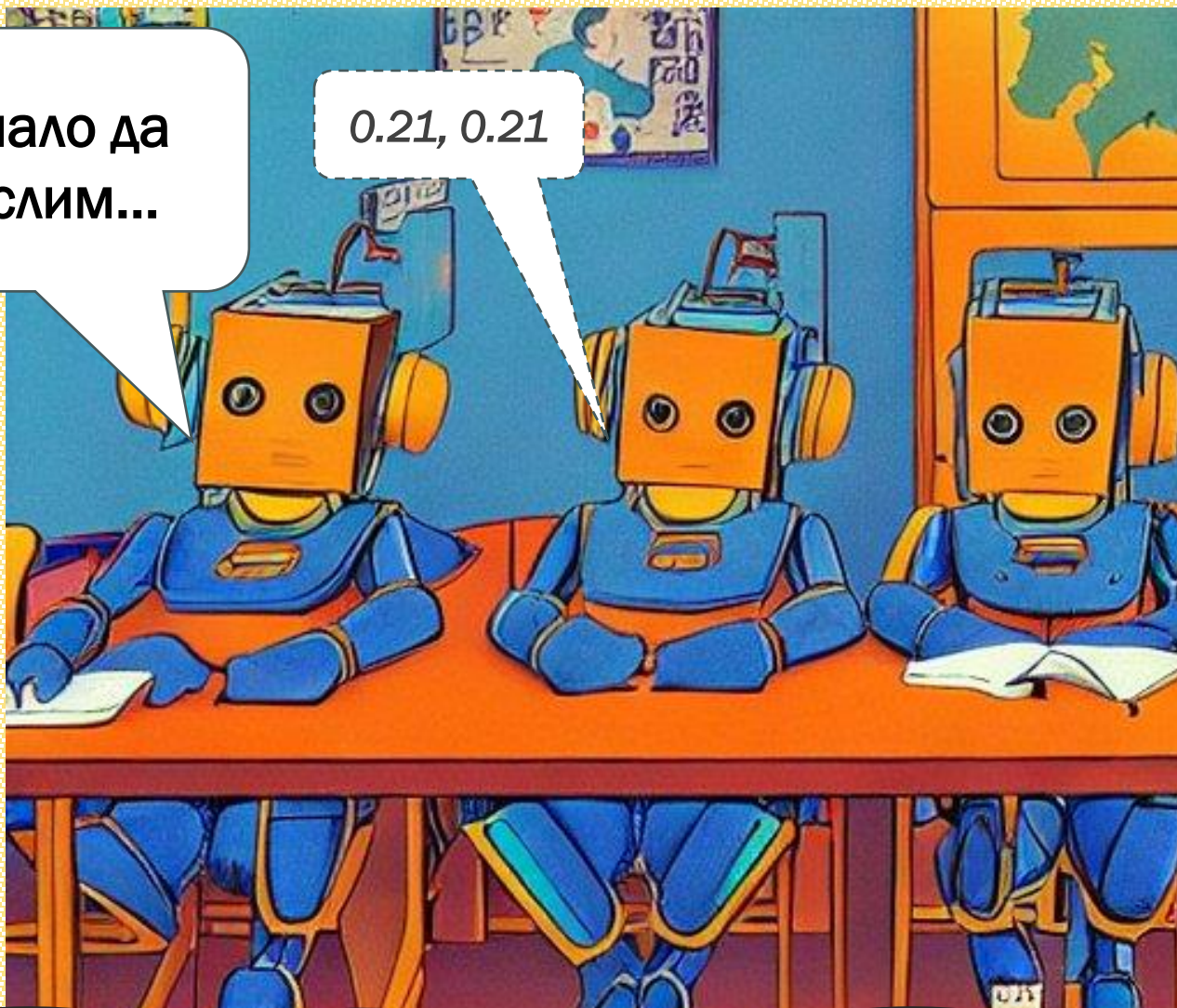
ХЕА-12, колика је
вероватноћа да је ова
реченица исправна?

Наставница пита
колика је вероватноћа
исправности реченице..

**ЧАС СРПСКОГ
-будућност-**

Само мало да
размислим...

0.21, 0.21

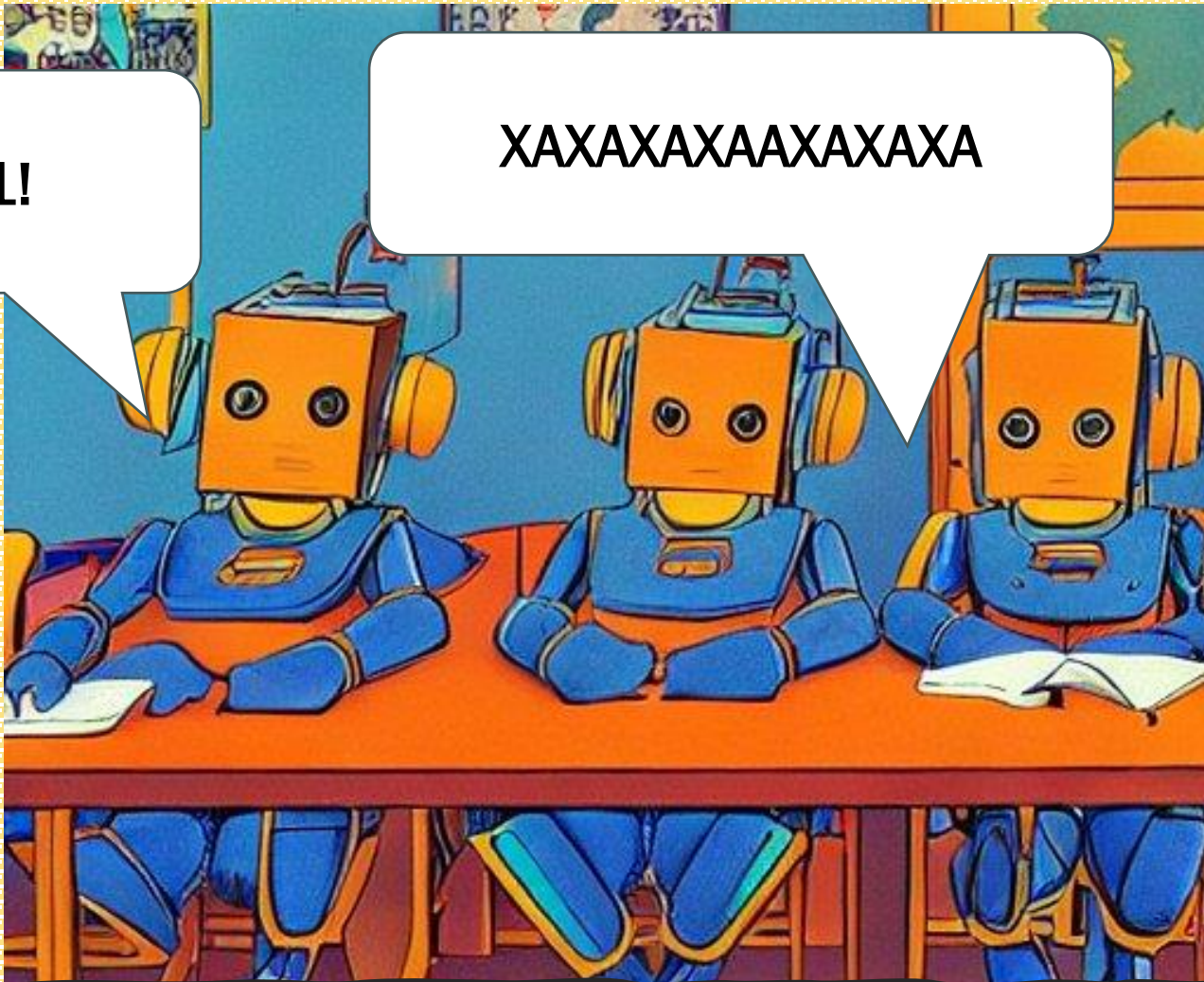


Од *деце* се очекује
одговор на скали од 0 до 1.

Што је много већи спектар
од буловског одговора.

0.21!

XAXAXAXAXAXAXA



Језички модели додељују нискама текста вероватноће на основу корпуса за обучавање, што не одговара нужно расподели вероватноћа у природном језику.

Оцењивање из српског у будућности??

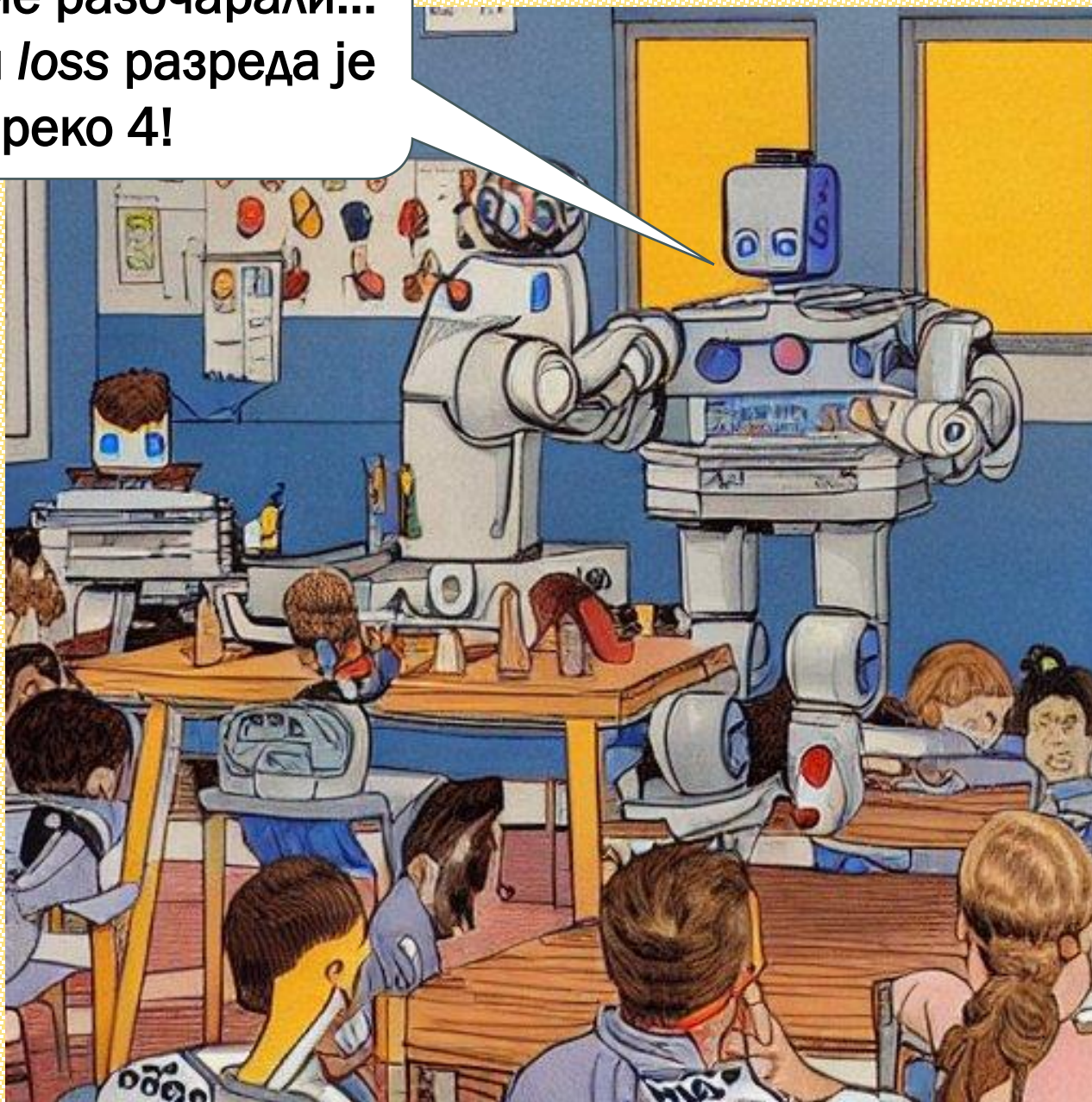
x – вектор "исправних" вероватноћа реченица

y – вектор вероватноћа која је ученик исписао на тесту

$$loss(x, y) = - \sum_{i=1}^n x_i \log y_i$$

$loss < 2$ је, на пример, потребан за петицу.

Јако сте ме разочарали...
Просечни *loss* разреда је
преко 4!

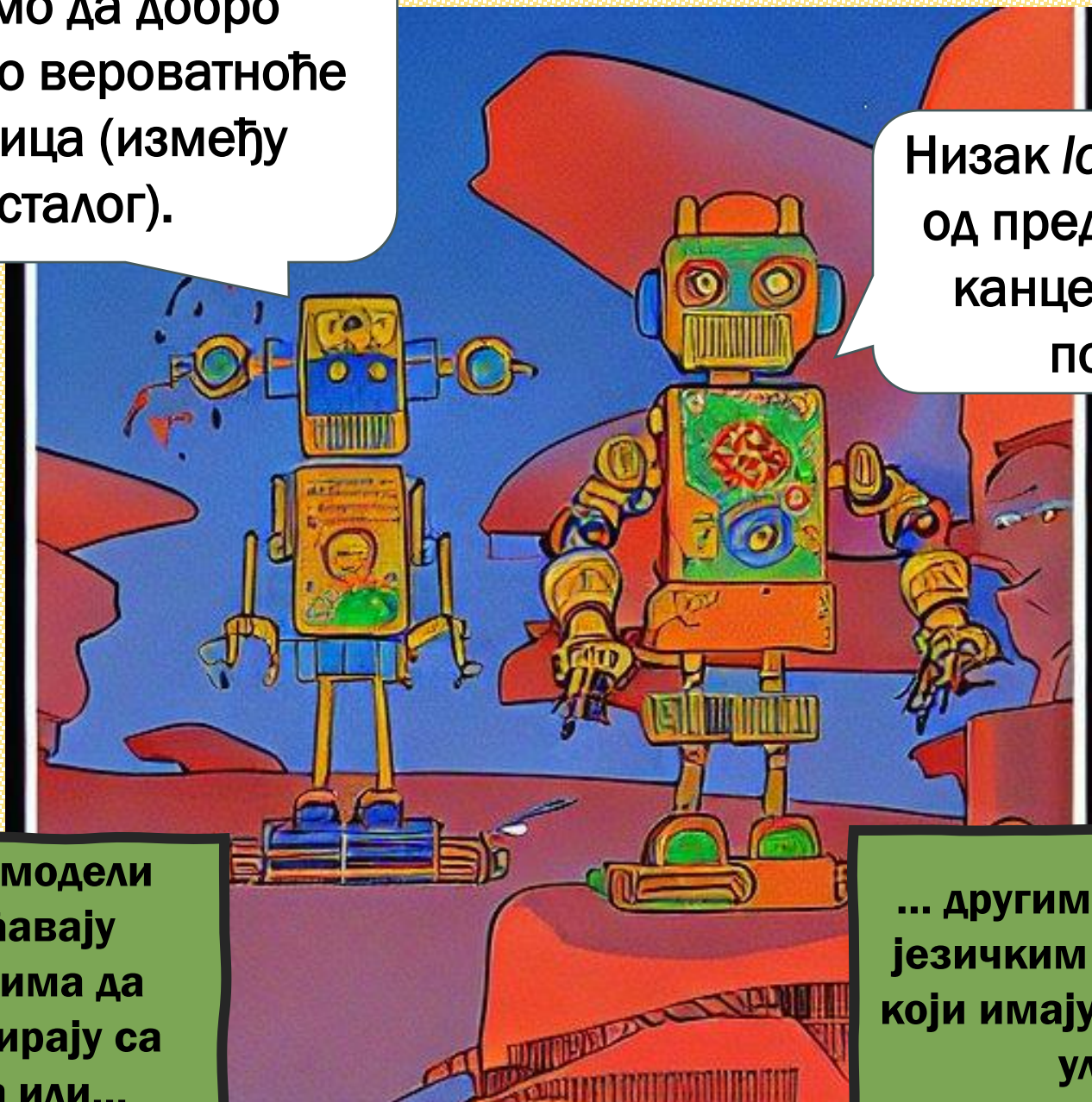


Желимо да добро
погађамо вероватноће
реченица (између
осталог).

Низак *loss* је један
од предуслова за
канцеларијски
посао.

Језички модели
омогућавају
рачунарима да
комуницирају са
људима или...

... другим модерним
језичким моделима,
који имају текстуални
улаз.

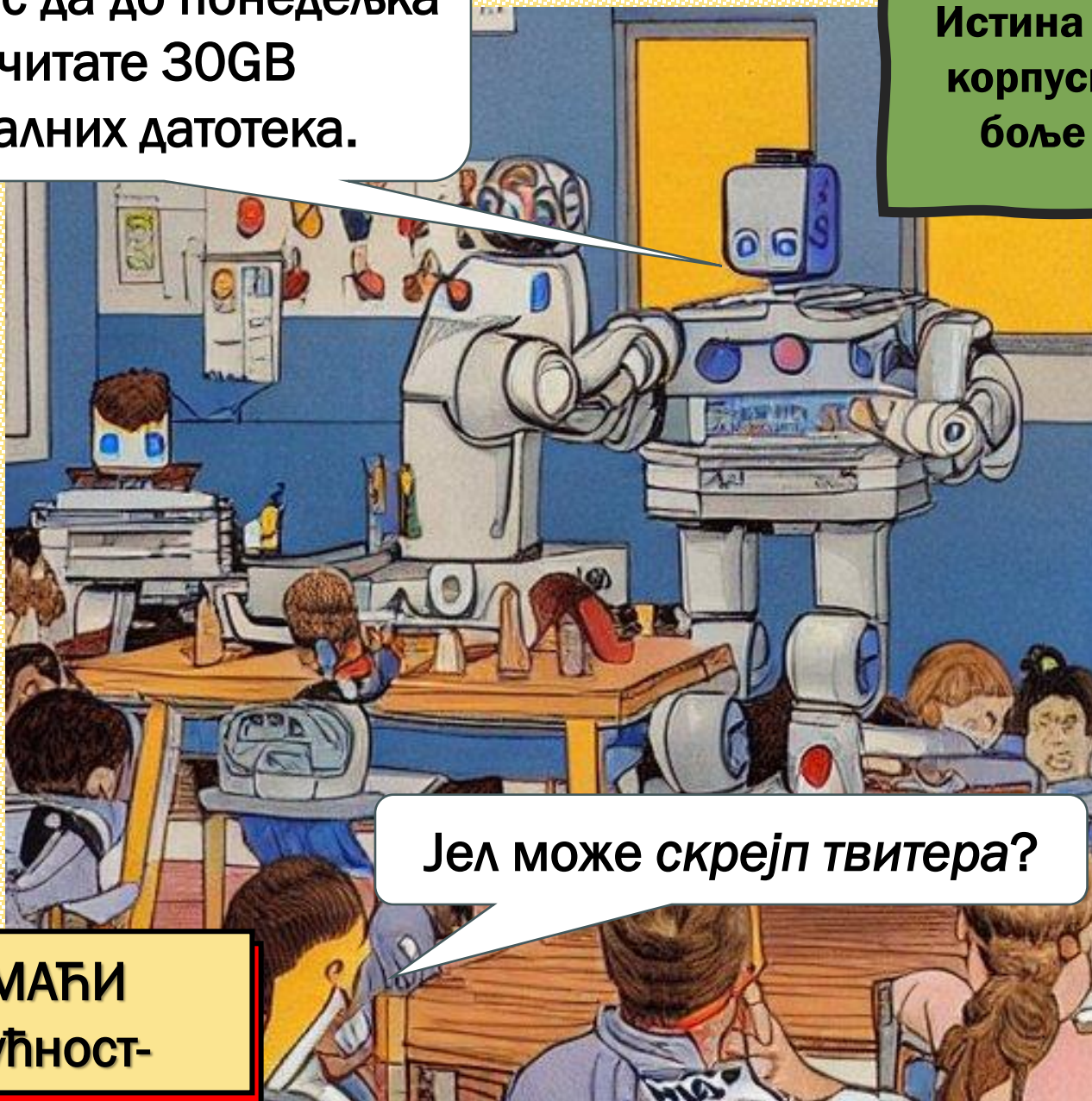


Шта је потребно за обучавање језичких модела?

- 1. Корпус квалитетног текста
(у што већем квантитету)**
- 2. Рачунарски ресурси
(што моћнији и што непрекиднији)**
- 3. Софтвер?**

Молим вас да до понедељка
прочитате 30GB
текстуалних датотека.

Истина је да већи
корпуси узрокују
боље моделе.

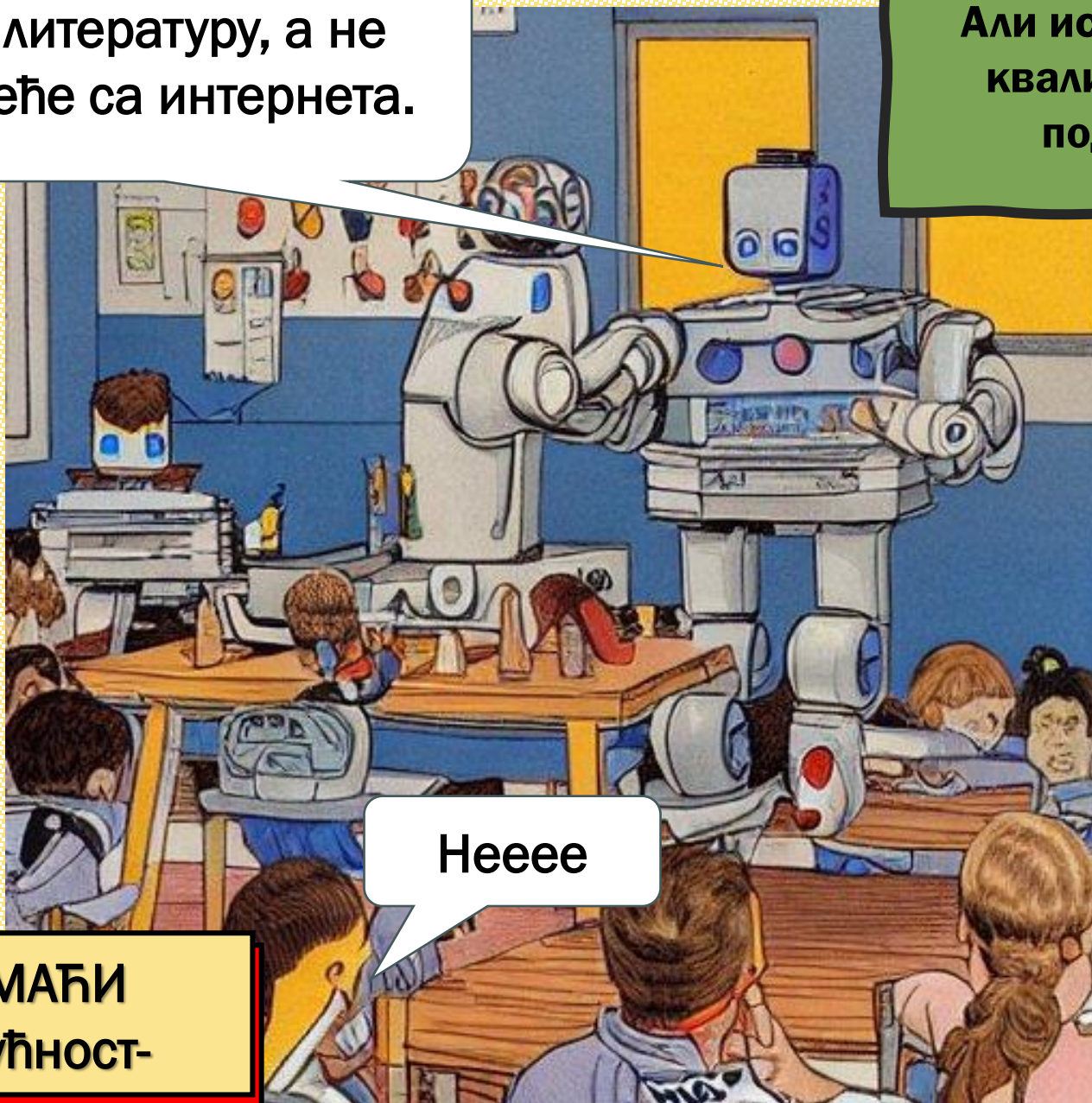


Јел може скрејп твитера?

ДОМАЋИ
-будућност-

Читајте литературу, а не само смеће са интернета.

Али исто тако и квалитетнији подаци.



Heeee

ДОМАЋИ
-будућност-



**Спор је,
али је наш.**

**Морамо му купити
нову графичку
картицу, заостаје за
разредом.**

**За обучавање се може
користити било који
рачунар, под условом да
модел и скуп за обучавање
могу да стану у активну
меморију.**

**Ипак, добра графичка
картица ће омогућити да се
то оствари много брже
(уколико се испуне
одређени услови).**

Препорука за софтвер?

1. [Трансформерска архитектура](#)
(уколико желите најбоље резултате);
2. Пајтон библиотека [transformers](#)
(лична препорука за обучавање и коришћење трансформерских модела);
3. [Scratch2LM](#), мој пројекат на *GitHub*-у за обучавање и дообучавање језичких модела
(уколико не желите да кренете од почетка).

Трансформерска архитектура?

1. Тренутно *state-of-the-art* моделирања језика (сви најбољи модели су засновани на њој);
2. Комплексна неуронска мрежа са милионима или чак милијардама параметара (углавном, више параметара—бољи резултати);
3. Заснива се на енкодеру (који векторизује улазни текст) и декодеру (који генерише нови, излазни текст)
(АЛИ НЕ НЕОПХОДНО. Многобројни модели засновани су само на енкодеру, или само на декодеру).

ЕНКОДЕР МОДЕЛИ

излаз је вектор

2.6123623847961426,
1.3682820796966553,
8.226987838745117,
0.8015447854995728,
1.9746443033218384,
-0.1605355143547058,
-1.6255472898483276,
-0.6387592554092407,
...

...

Сине, како
је било у
школи?

Енкодер за сваку ниску текста враћа одговарајући вектор, који је представља.



Енкодер модели?

- На излазу производе векторе;
- Обучавају се (обично) над листом реченица;
- Дообучавају се (засновано на излазним векторима) на задацима обележавања речи, класификације реченица, али и допуњавања недостајућег текста;
- Најпознатији примерци: BERT, RoBERTA, ELECTRA...
- Најбољи модели за српски на платформи *huggingface*: [xlm-roberta](#)? [BERTić](#)? [SRoBERT-a-F](#)? (сви наведени модели су мање или више вишејезични)*

*Ускоро се очекује објављивање модела [Bertović-base](#) и [Bertović-large](#), обучаваних специјално за српски језик.

ДЕКОДЕР МОДЕЛИ

излаз је наставак улаза

Упитао је са болом у очима. Тог тренутка је знао да...

Сине, како је било у школи?



Декодер модели за улазну ниску текста враћају један или више вероватних наставака.

Декодер модели?

- Обучавају се над листом n -грама токена (обично текст исцепкан на једнаке делове);
- На излазу производе наставак текста;
- Дообучавају се за генерисање специфичних стилова (рачунарски код, хаику песме...);
- Најпознатији примерци: GPT(2,3,4...), LLaMA, PALM (ChatGPT и BARD врло вероватно нису декодер модели)
- Најбољи модели за српски: [sr-gpt2](#), [gpt2-sr-lat](#) и [sr-gpt2-large](#), највећи декодер модел за српски са преко 700 милиона параметара.

ЕНКОДЕР-ДЕКОДЕР МОДЕЛИ

излаз је трансформација

Mon fils, comment
était l'école?

Сине, како
је било у
школи?

Најкомплекснији модели.

Генеришу излаз на основу улаза.



ЕНКОДЕР-ДЕКОДЕР МОДЕЛИ

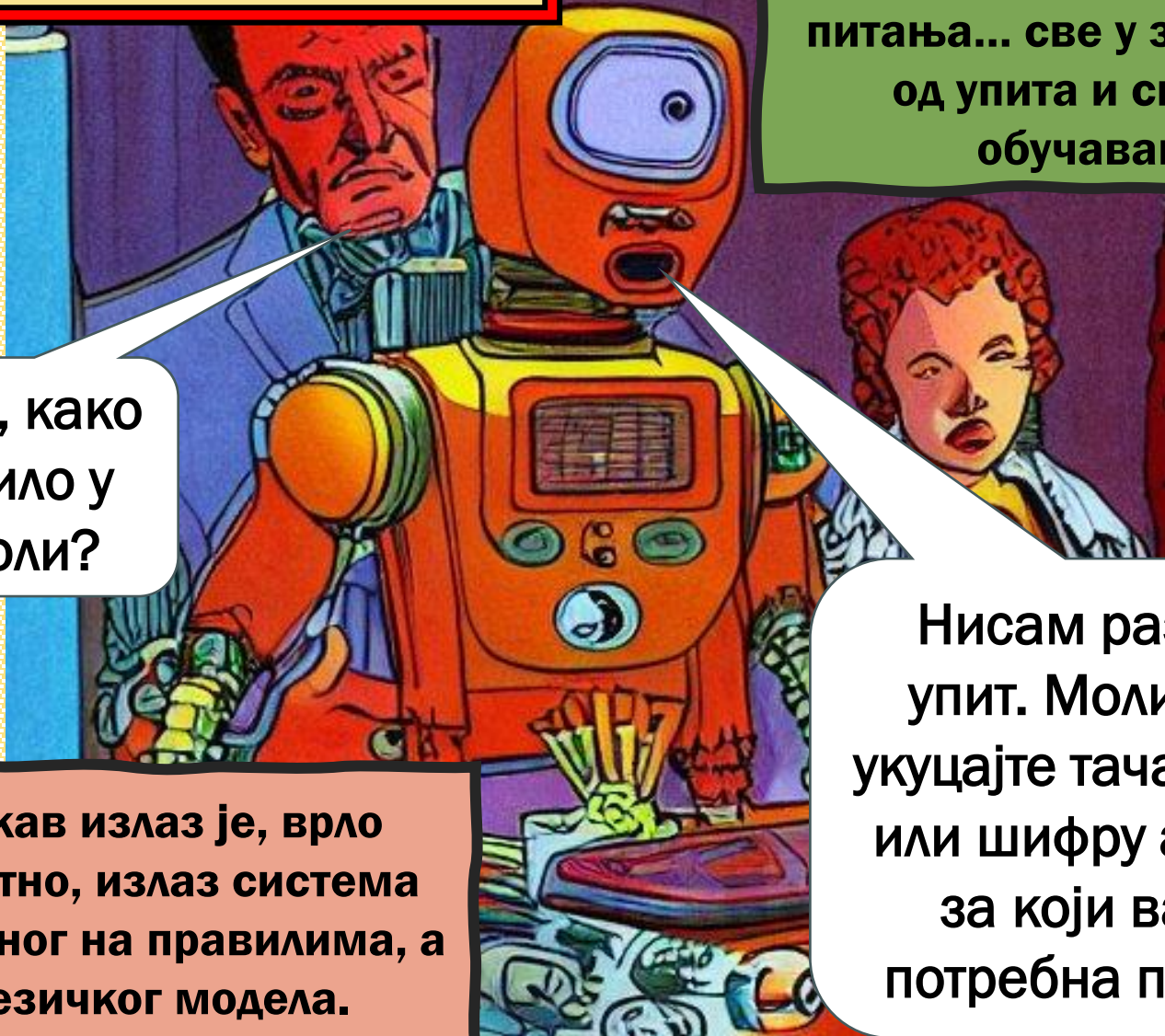
излаз је трансформација

Излаз могу бити трансформације текста, преводи, одговори на питања... све у зависности од упита и скупа за обучавање.

Сине, како је било у школи?

Овакав излаз је, врло вероватно, излаз система заснованог на правилима, а не језичког модела.

Нисам разумео упит. Молим вас укуцајте тачан назив или шифру артикла за који вам је потребна помоћ.*



Енкодер-Декодер модели?

- **Обучавају се над паровима реченица;**
- **На излазу производе нови текст;**
- **Користе се за:**
 - **машинско превођење**
 - **одговарање на питања**
 - **сажимање текста**
 - **трансфер стила...**
- **Најпознатији примерци: BART, T5, Galactica (ChatGPT и BARD су вероватно овог типа)**
- **За српски тренутно нема познатих модела...**

МОДЕЛИ ЗА СРПСКИ НАСПРАМ ВИШЕЈЕЗИЧНИХ МОДЕЛА

	bertovic	srberta	maced.	sroberta	sroberta-x	xlm
bertovic-base	0.729	0.591	0.610	0.628	0.705	0.572
JelenaTosic/SRBerta	0.408	0.644	0.496	0.377	0.408	0.426
maced./sr-roberta-base	0.318	0.387	0.454	0.302	0.317	0.326
Andrija/SRoBERTa	0.279	0.219	0.222	0.275	0.276	0.246
Andrija/SRoBERTa-base	0.410	0.352	0.358	0.378	0.429	0.365
Andrija/SRoBERTa-L	0.573	0.486	0.500	0.534	0.629	0.485
Andrija/SRoBERTa-F	0.647	0.576	0.592	0.593	0.702	0.544
Andrija/SRoBERTa-XL	0.609	0.535	0.558	0.559	0.660	0.518
<i>xlm-roberta-base</i>	0.465	0.395	0.407	0.470	0.462	0.503
<i>xlm-roberta-large</i>	0.515	0.436	0.453	0.513	0.507	0.545

Тачност различитих модела (уз различите токенизаторе) на задатку допуњавања недостајућег текста у роману Дечко.

ВИШЕЈЕЗИЧНИ МОДЕЛ

Један од разлога за
обучавање модела за
српски језик.

Да, говорим
преко сто
светских
језика!

Не, не знам
одговор на
твоје питање!



Моје обучавање језичких модела за српски

Састав корпуса:

1. СрпКор2013 (50%)
2. СрпКор2021 (20%)
3. ВикиКорпус (28%)
4. СрпЕЛТеК (2%)

**Корпуси
квалитетног
текста**

Рачунарски ресурси:

Серверске машине са
Nvidia графичким картицама:

- RTX 2060 6gb
- RTX 3060 12gb
- RTX 4090 24gb

**Постепено
унапређење
ресурса**

Обучени модели:

gpt2-srlat – основни GPT2 модел

gpt2-srlat-sem – семантички модел

gpt2-srlat-synt – синтаксички модел

+ bertović-base (енкодер модел
обучен на истом скупу података)

**Обучени различити
модели над
различитим
репрезентацијама
корпуса.**

ОСНОВНИ ГПТ МОДЕЛ

Па ви сте
робот,
господине!

Ова ниска текста ме
не изненађује.
перплексност = 13.7

Трансформерски
модели одашиљају
њихову изненађеност
улазном ниском
(перплексност).

узвикну он.

Контролни примерак, предефинисана подешавања.



СЕМАНТИЧКИ ГПТ МОДЕЛ

дугме река
одело брати

дообучавање основног модела
непроменљиве речи уклоњене
променљиве замењене лемама

Ово нема никаквог
смисла.
перплексност = 492

Перплексност је
реципрочна
вероватноћа.

пар нојев перје
насеље

Циљ је било моделирање
семантичког аспекта.



СИНТАКСИЧКИ ГПТ МОДЕЛ

ADV Yys од

Разумно.
перплексност = 54

Брзо узми од
човека. ?

ms2v.

дообучавање основног модела
променљиве речи замењене су
граматичким категоријама (шифрама)

Циљ је било моделирање
граматичности текста.



Мудрост гомиле и композитни језички модели

Композитни језички модели се заснивају на претпоставци да више глава боље размишљају од једне.



Један од популарнијих експеримената је давање групи људи да погоде колико слаткиша има у тегли.

Погодите колико има слаткиша у тегли

МУДРОСТ ГОМИЛЕ



100

18

84

141

МУДРОСТ ГОМИЛЕ

Различите индивидуе дају различите одговоре, на основу сопственог предосећаја, логике или искуства.



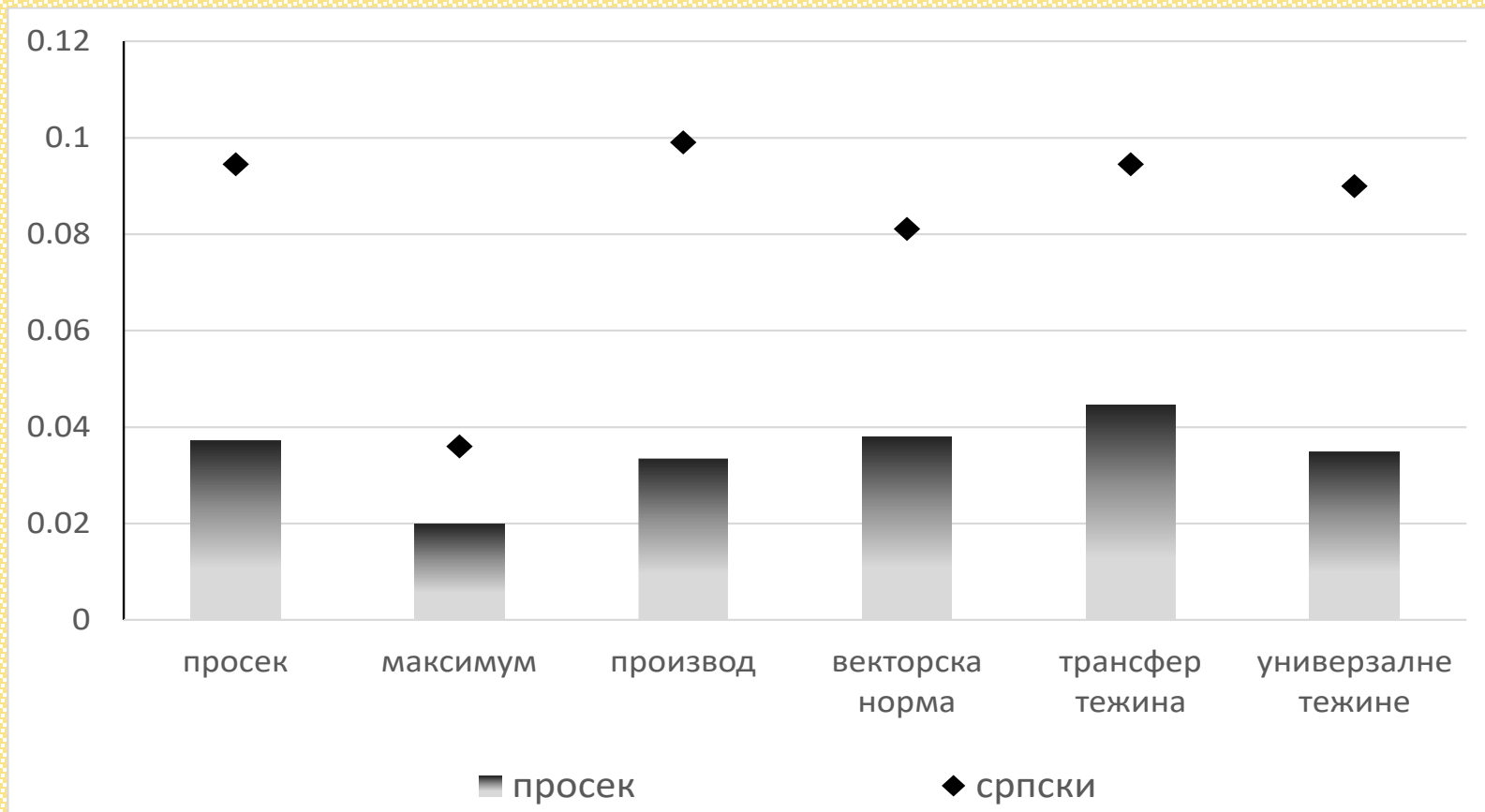
Заправо, тачан одговор

$$\text{је } \approx \frac{1}{n} \sum_{i=1}^n x_i$$

Тачан резултат је често приближан аритметичкој средини одговора из гомиле.

МУДРОСТ ГОМИЛЕ

Ово је показано као делотворно у моделовању језика (и то поготово српског) на пробном експерименту.



Проценти унапређења добијеног коришћењем композитне архитектуре за седам европских језика, укључујући и српски.*

*Škorić, Mihailo, et al. "Parallel stylometric document embeddings with deep learning based language models in literary authorship attribution." Mathematics 10.5 (2022): 838.

Колика је вероватноћа
да је ова реченица
исправна?

Различити модели
дају различите мере
вероватноће за исту
нисуку текста.



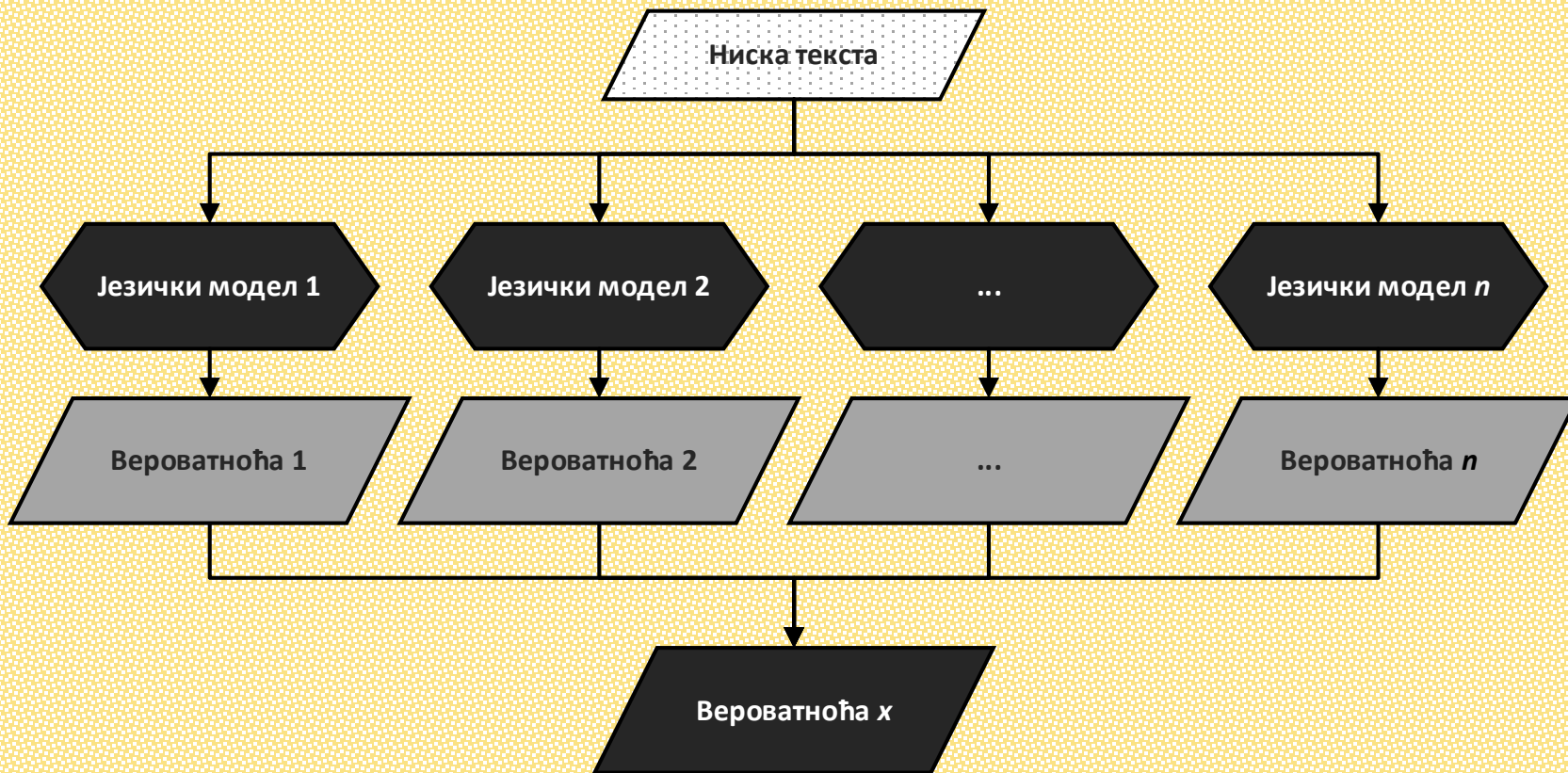
0.47

0.76

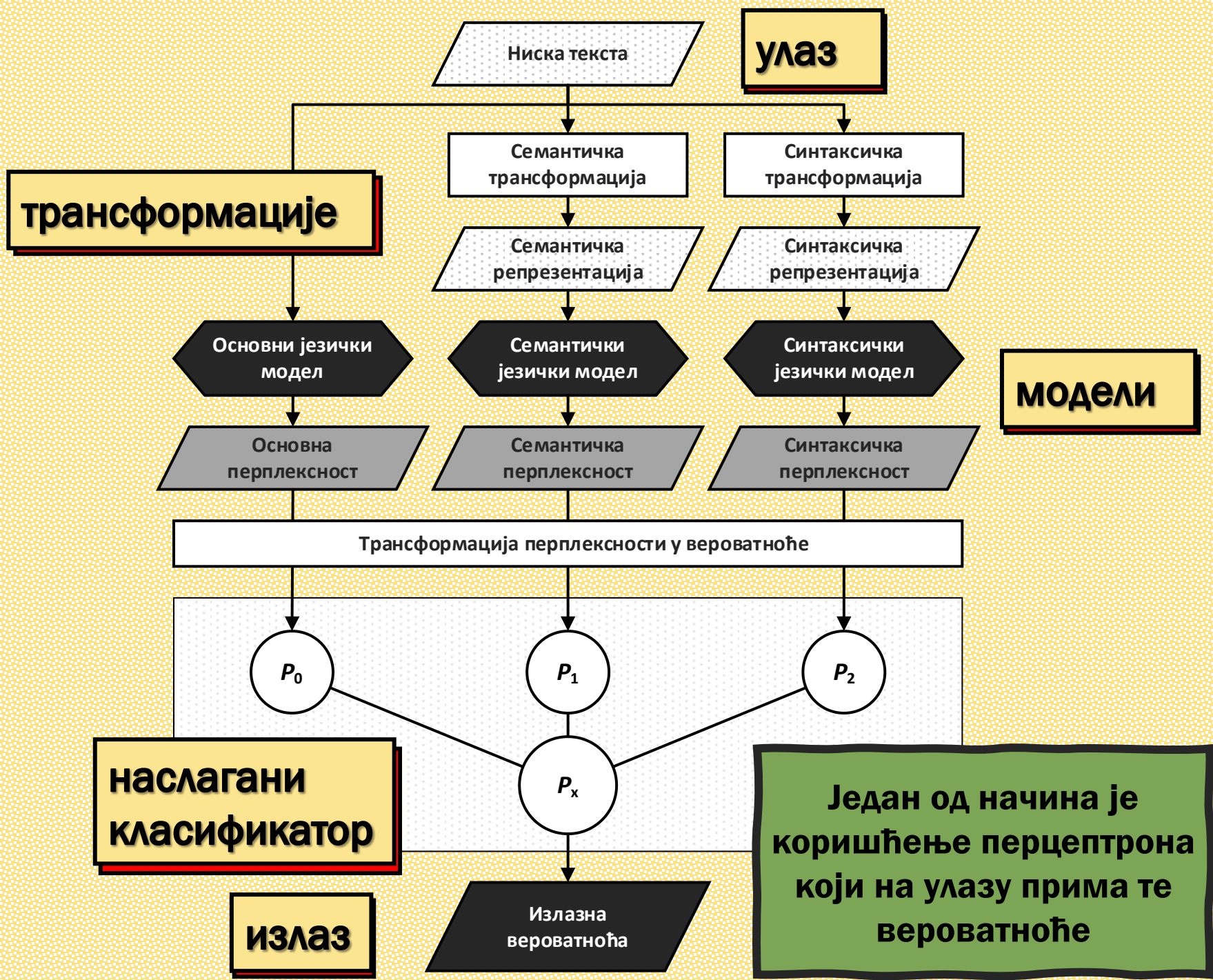
0.32

комполитна архитектура,
паралелно увезаних језичких модела

Композитна архитектура, паралелно увезаних језичких модела

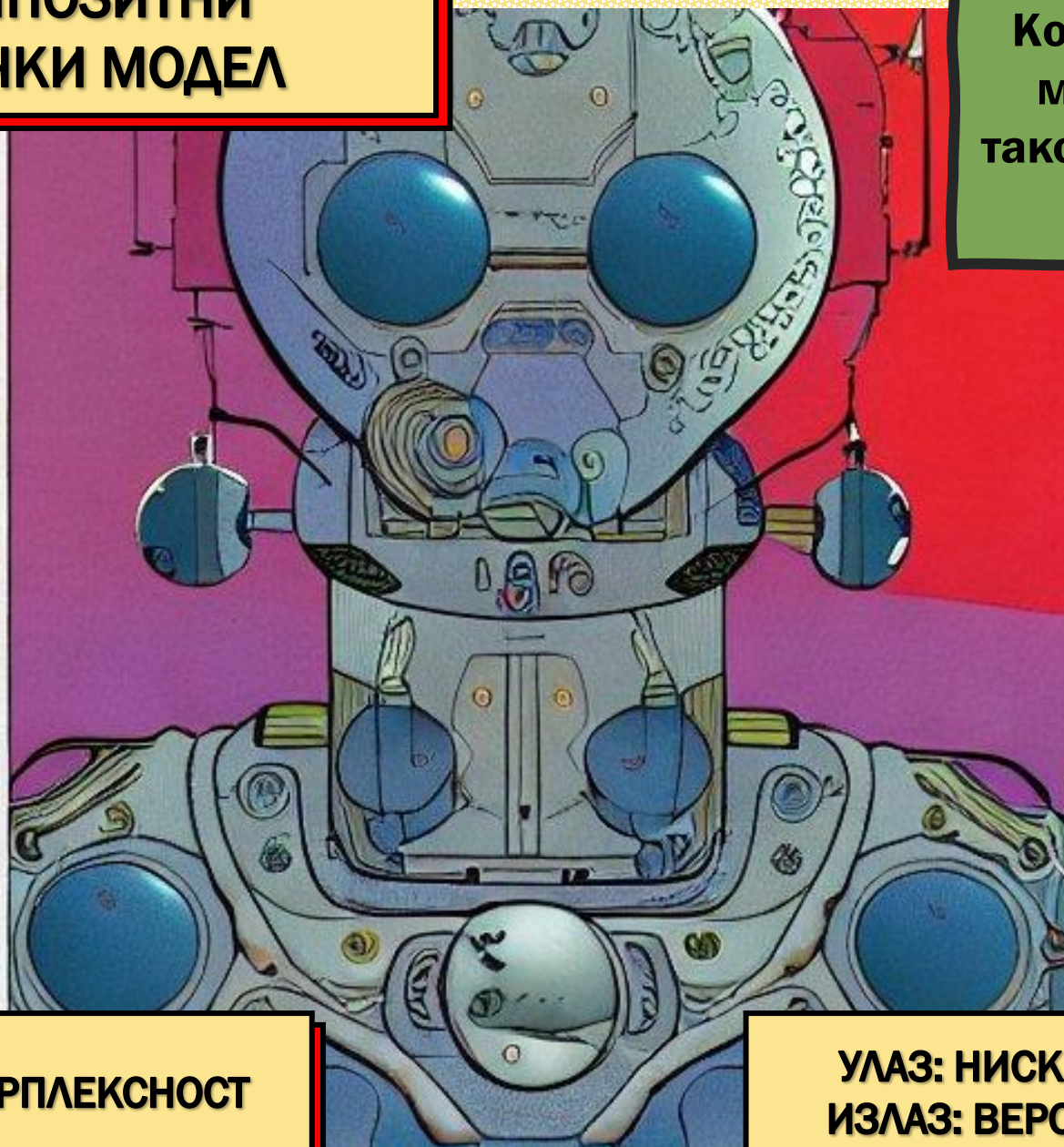


А онда се те вероватноће на неки начин збрајају.



КОМПОЗИТНИ ЈЕЗИЧКИ МОДЕЛ

Композитни
модели су
такође језички
модели



МЕДИЈ: ПЕРПЛЕКСНОСТ

УЛАЗ: НИСКА ТЕКСТА
ИЗЛАЗ: ВЕРОВАТНОЋА

**Унапређења која композитни модел даје у односу на основни,
на различитим задацима.**

затак	процент повећања тачности	процент смањења грешке
детекција синтаксички неисправних реченица (према облицима речи)	3.15	18.00
детекција синтаксички неисправних реченица (према редоследу речи)	0.46	3.67
детекција семантички неисправних реченица	0.21	0.86
разликовање синтаксички и семантички неисправних реченица (први тест)	10.10	23.61
разликовање синтаксички и семантички неисправних реченица(други тест)	9.98	28.22
<i>детекција неисправних реченица уопштено</i>	/	/
<i>разликовање експертских и машинских превода</i>	6.00	6.44

На последња два задатка модел не конвергира.

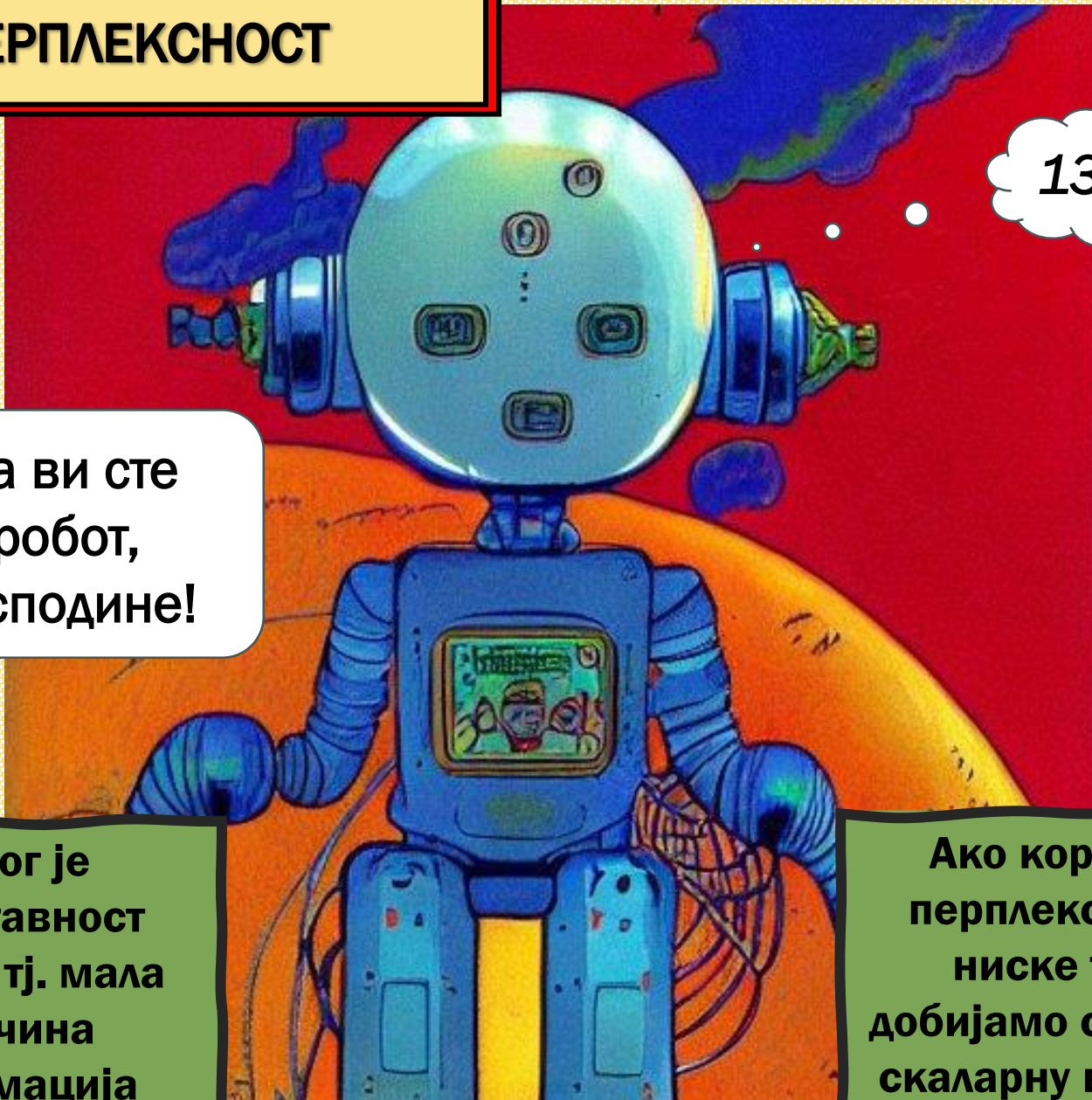
ПЕРПЛЕКСНОСТ

13.7

Па ви сте
робот,
господине!

Разлог је
једноставност
модела, тј. мала
количина
информација

Ако користимо
перплексност, од
ниске текста
добивамо само једну
скаларну вредност..



ВЕКТОР ПЕРПЛЕКСНОСТИ

Па ви сте

ви сте робот

сте робот,

робот,
господине

, господине!

Количину информација можемо повећати бомбардовањем модела н-грамима текста, уместо једне ниске.

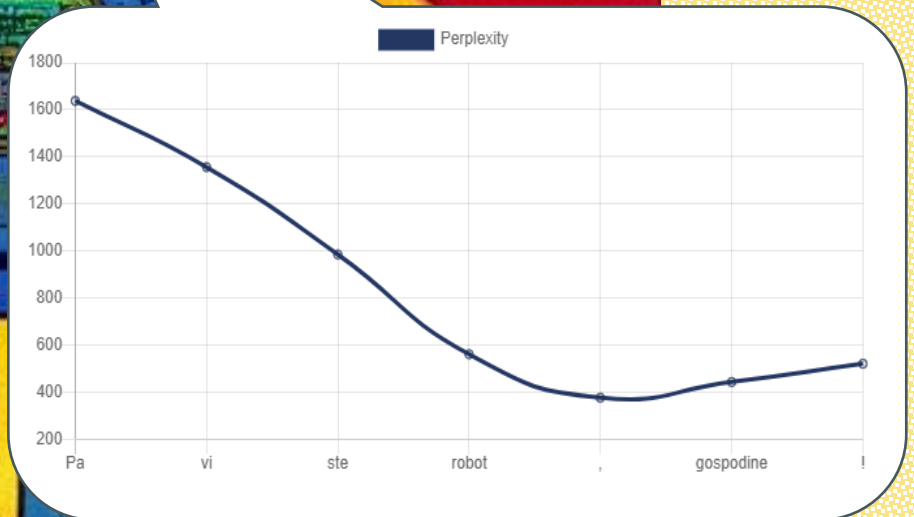
44.4

50.4

33.0

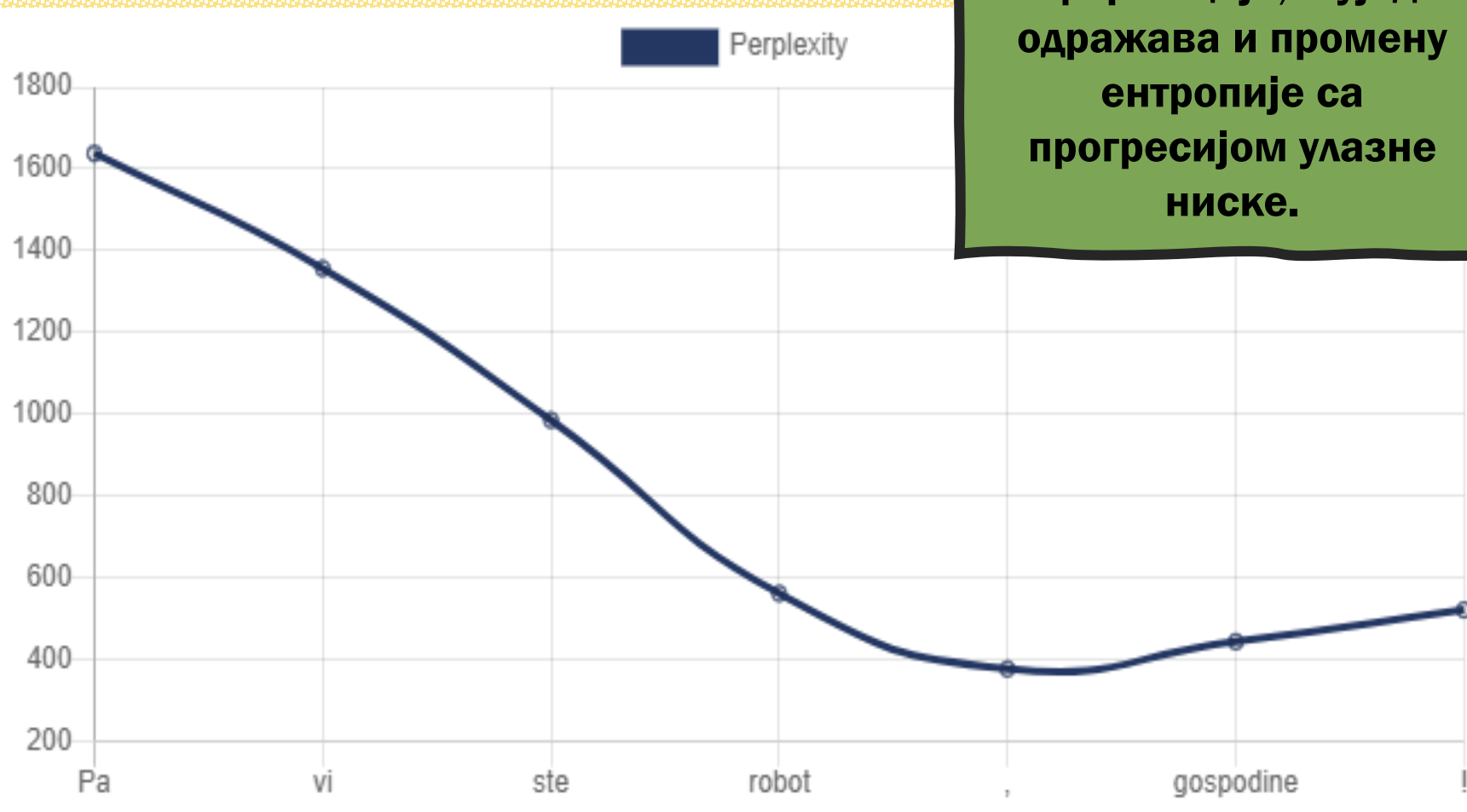
53.3

35.6



ВЕКТОР ПЕРПЛЕКСНОСТИ

Добијени вектор перплексности носи већу количину информација, а уједно одражава и промену ентропије са прогресијом улазне ниске.



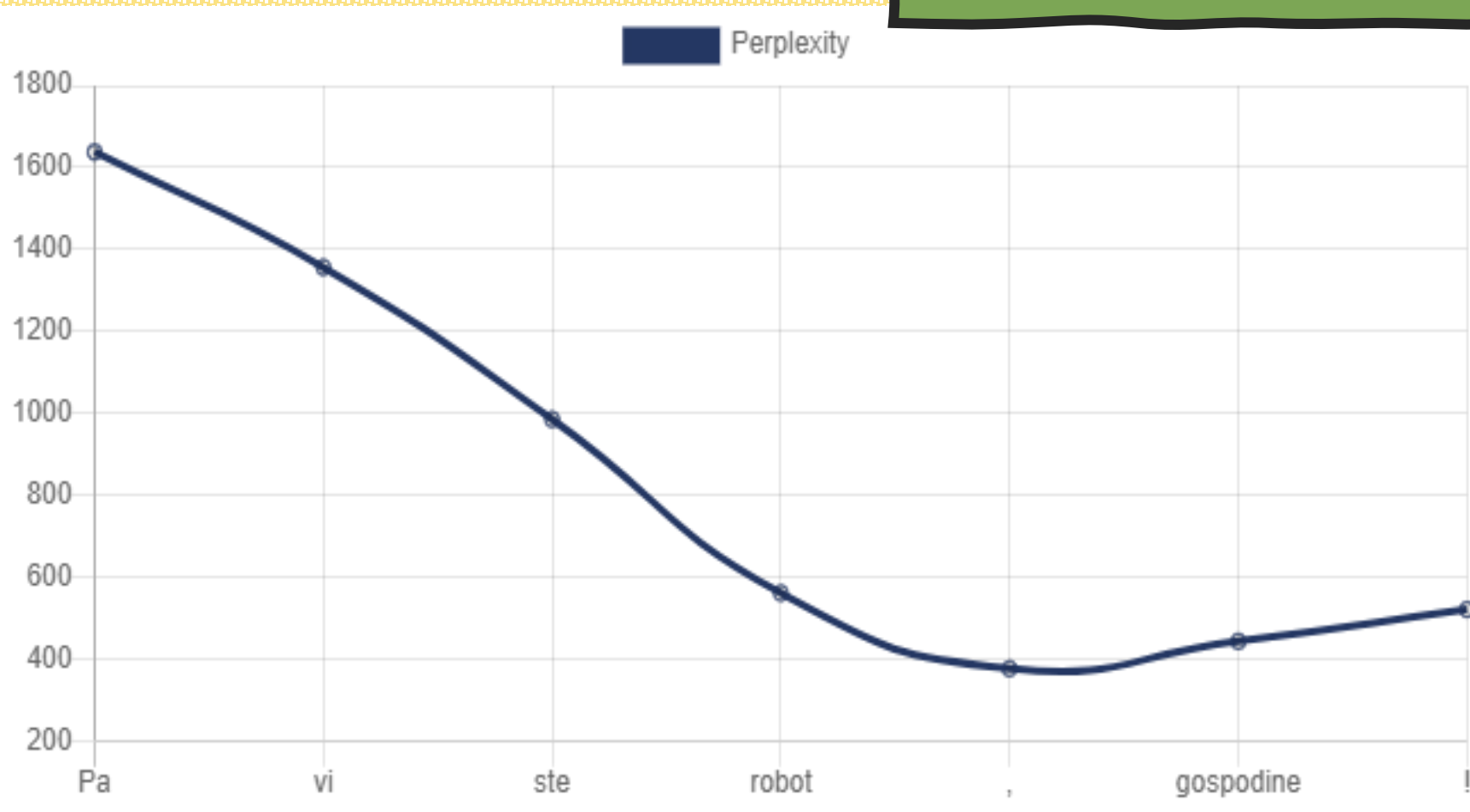
**Што омогућава
детекцију потенцијалних
грешака у тексту.**

затак	насумично	Основни модел	Семантички модел	Синтаксички модел
детекција избрисане речи	0.058	0.104	0.053	0.076
детекција уметнуте речи	0.031	0.173	0.036	0.068
детекција замењене речи	0.020	0.186	0.030	0.060

**При чему се тачност на
овом задатку повећава
и до 9 пута у односу на
насумични одабир.**

ВЕКТОР ПЕРПЛЕКСНОСТИ

Вектори омогућавају
обучавање комплекснијих
модела попут коволуционих
неуронских мрежа.



улаз

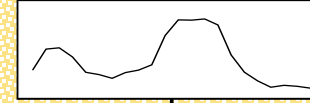
Ниска текста

вектори перплексности

Основни ВП
(\overline{PP}_0)

Семантички ВП
(\overline{PP}_1)

Синтаксички ВП
(\overline{PP}_2)



Трансформација ВП у векторе вероватноће

Уједначавање величина и комбиновање вектора вероватноће



Конволуција (4*3)

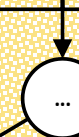
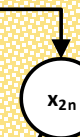
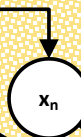
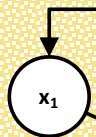
Конволуција (5*3)

Конволуција (6*3)

Сажимање

Сажимање

Сажимање



излаз

Излазна вероватноћа

Вектори омогућавају обучавање комплекснијих модела попут рекуретних или коволуционих неуронских мрежа.

КОМПОЗИТНИ
вектор

КОНВОЛУЦИЈА

наслaгани
класификатор

КОМПОЗИТНИ ЈЕЗИЧКИ МОДЕЛ 2.0

Чиме се добија
комплекснија
структура са већом
моћи учења.



МЕДИЈ: ВЕКТОРИ
ПЕРПЛЕКСНОСТИ

УЛАЗ: НИСКА ТЕКСТА
ИЗЛАЗ: ВЕРОВАТНОЋА

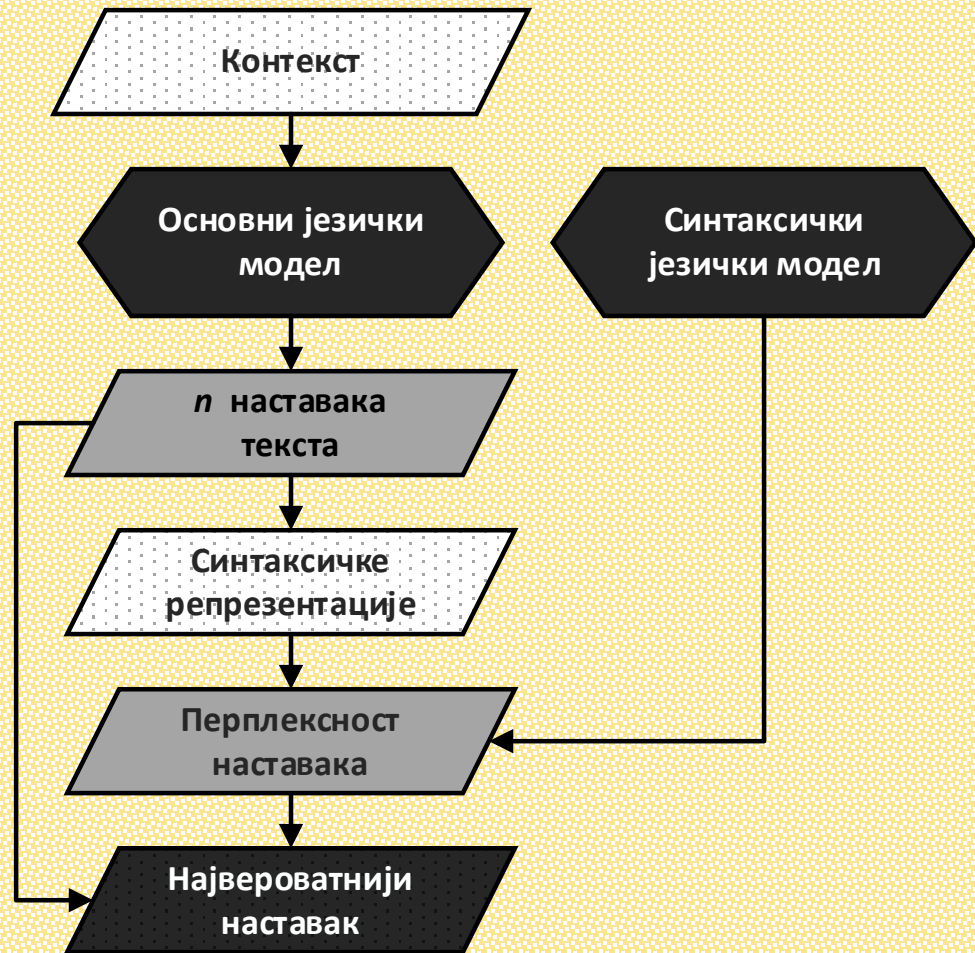
**Која може да савлада
комплексније задатке.**

затак	процент повећања тачности	процент смањења грешке
детекција неисправних реченица уопштено	9.89	48.24
разликовање експертских и машинских превода	8.74	9.67

КОМПОЗИТНО ГЕНЕРИСАЊЕ ТЕКСТА

Што се тиче генерисања текста коришћењем више обучених модела за српски...

...најбоље резултате постиже комбинација основног и синтаксичког модела.



Основни модел
генерише n
наставака
текста...

Нисам сигуран који је
од ових наставака
најбољи...

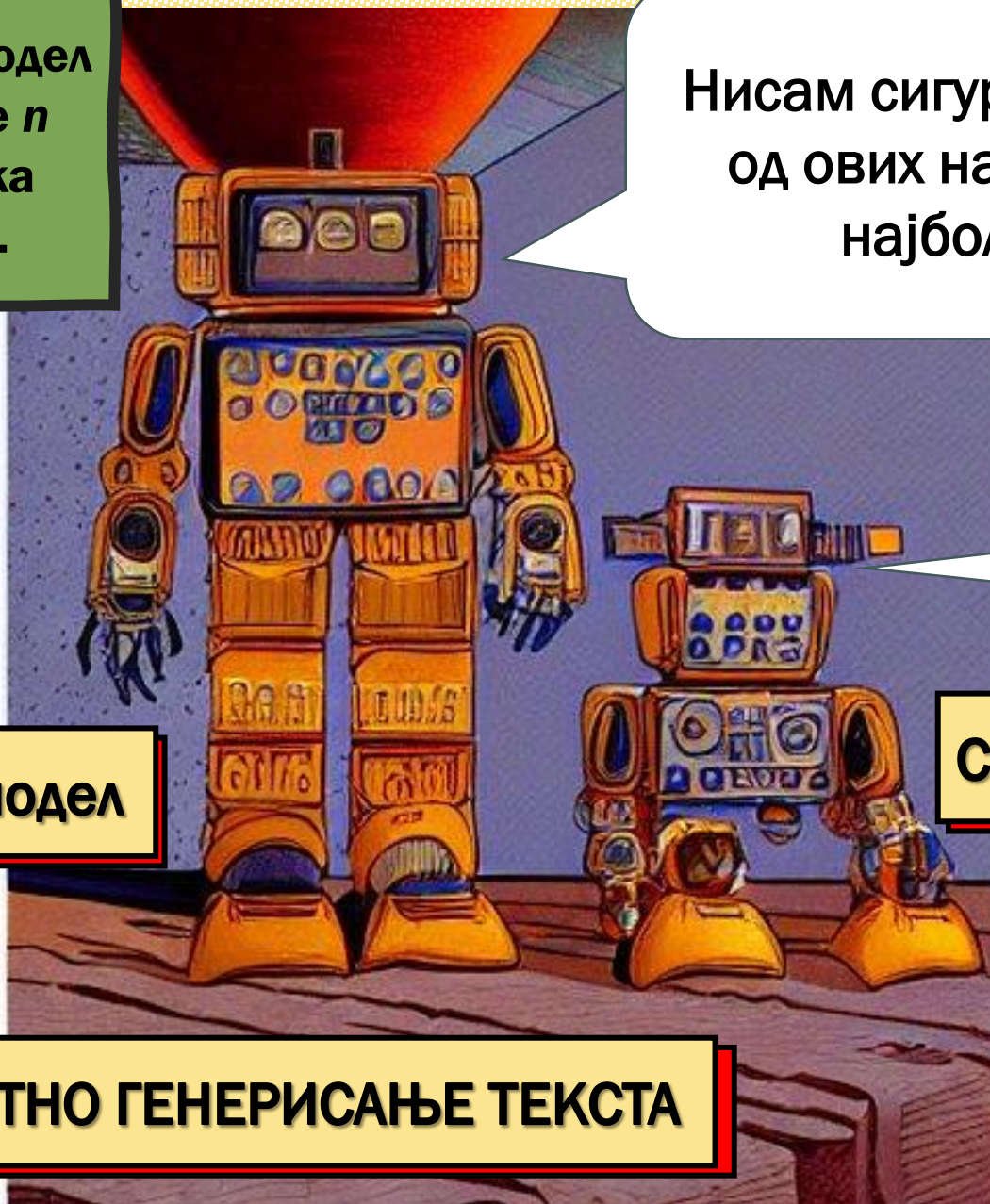
Други.

Основни модел

Синтаксички модел

КОМПОЗИТНО ГЕНЕРИСАЊЕ ТЕКСТА

...а синтаксички
каже који од
наставака је
најисправнији.



*Suma
sumarum...*

Било да је у питању
евалуација или генерисање
текста, када се удружимо,
постигемо боље резултате!

И вектори
перплексности
су кул!

И ТАКО...

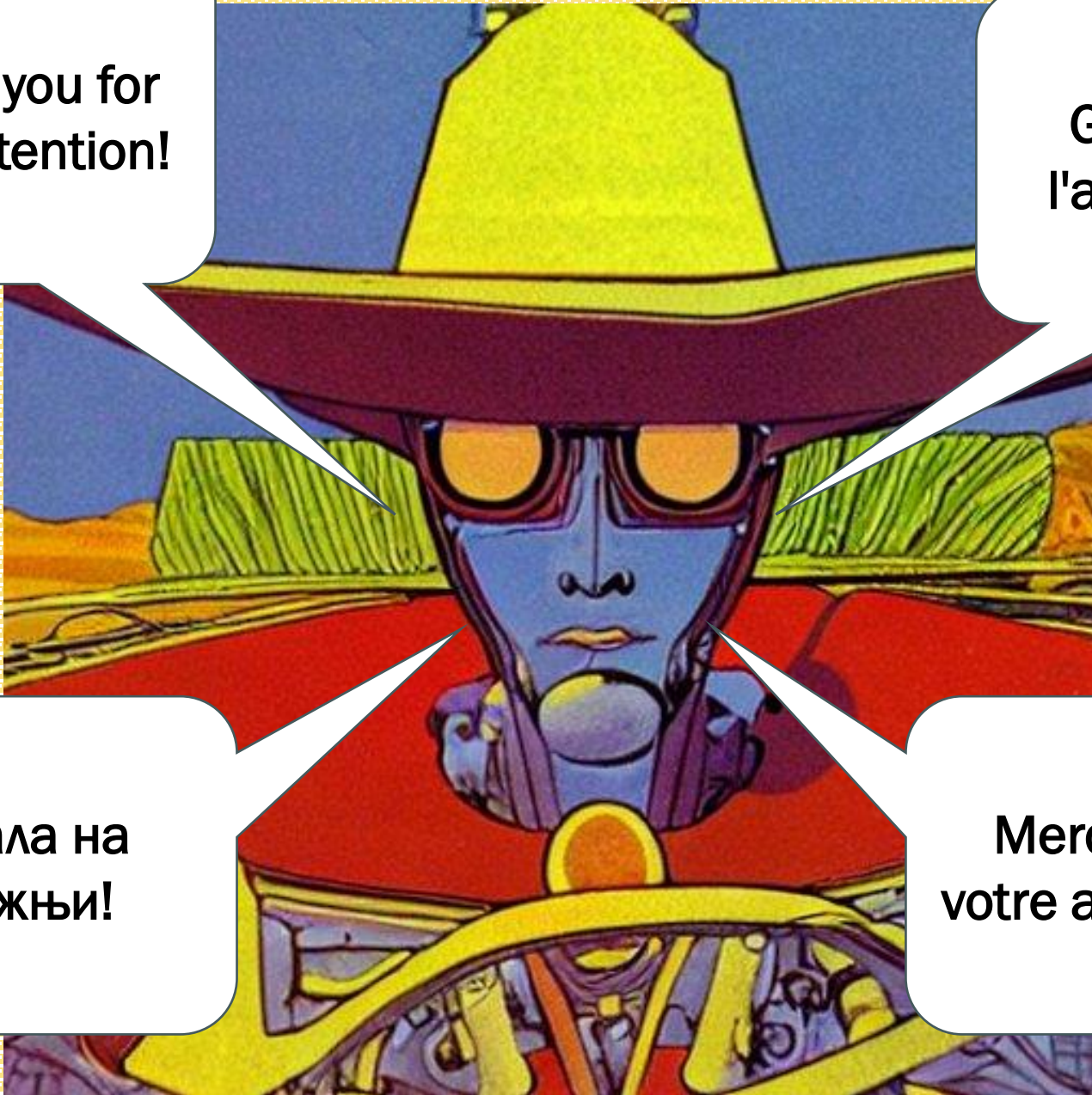


**Thank you for
your attention!**

**Grazie per
l'attenzione!**

**Хвала на
пажњи!**

**Merci pour
votre attention!**



PS:

МУЛТИМОДАЛНА ЕВАЛУАЦИЈА

