

MI Szakmai délután



OTP Bank
2023.09.14

Agenda

Érkezés 13:45-14:00

Előadás 14:00-15:30

Szünet 15:30-15:45

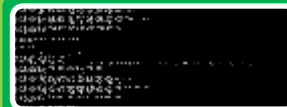
Kerekasztal beszélgetés 15:45-17:00



Bemutkozás: miért foglalkozik egy Bank a nyelvi modellekkel?
Előadó: **Kaszás Zoltán** -OTP Bank Ügyvezető Igazgató



OTP-s NLP „Szuperszámítógép” projekt céljai, tagjai
Előadó: **Nagy Zsombor**; Szenior projektvezető



Nyelvi modellek: miért szükséges a Magyar NLP modell?
Előadó: **OTP AI CoE**



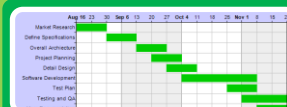
OTP szuperszámítógép választásának szempontjai
Előadó: **Nagy Zsombor**; **OTP AI CoE**



OTP szuperszámítógép építés lépései
Előadó: **Nagy Zsombor**



Alkalmazott biztonsági megoldások
Előadó: **Oláh István György** IT kockázatkezelési vezető Tanácsadó; Tanár



Kutatási terv és jelen állapot
Előadó: **Nagy Zsombor**



Mi várható eredményként
Előadó: **OTP AI CoE**

A kutatási projekt átfogó célja a digitális képességek fejlesztése a mesterséges intelligencia segítségével, a természetes és a **gépi nyelv közötti átjárás megteremtése és hogy létrejöjjön a mesterséges intelligencia-algoritmusok számára használható magyar nyelvi modell**. Ezzel a OTP csoport óriási lépést tesz a technológiában betöltött vezető szerep felé, egyúttal felkarolja a hazai mesterséges intelligencia-kutatás és -fejlesztés fellendítésének ügyét és elősegíti a magyar nyelvű mesterséges intelligencia alapú megoldások létrehozását is.

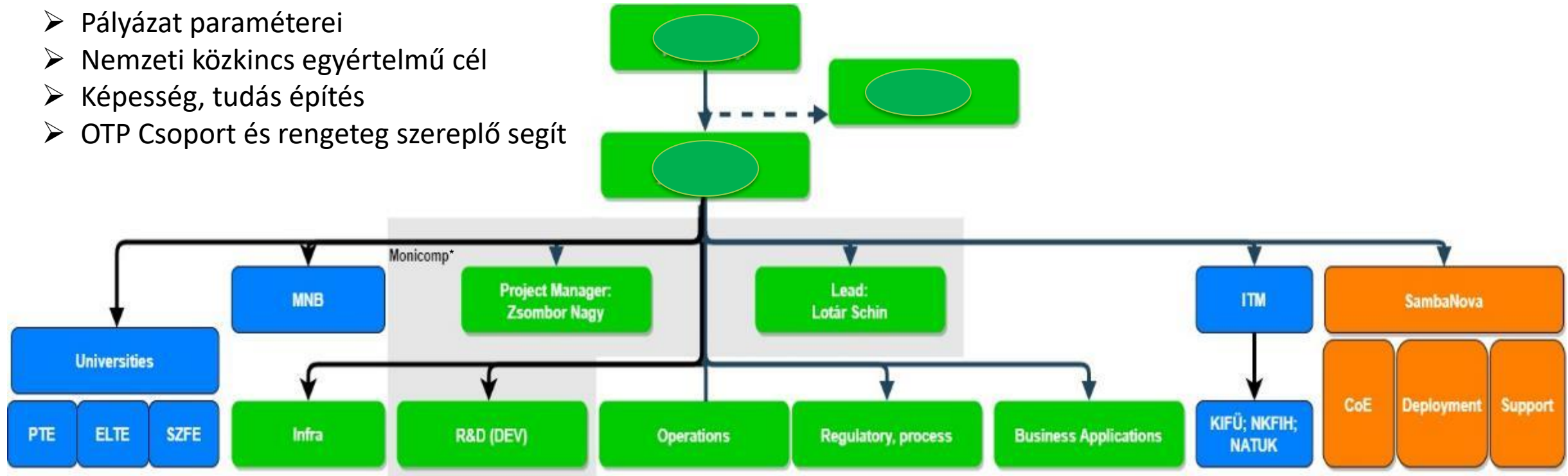
A magyar beszélt nyelv felhasználói köre relatíve szűk, ezért a globális piac nem érdekelt jó minőségű magyar nyelvfeldolgozás piaci alapú előállításában, emiatt a fejlesztés a globális kiszolgáltatottság mérséklésére irányul, mind a magyar nyelvfeldolgozás, mind az adatkezelés és a platform technológiák használata területén. A fejlesztés arra a hiányterületre épül, hogy az MI kutatásánál, fejlesztésénél a világnyelvekre, különösen az angolra fókuszálnak. Ez viszont az olyan kevés ember által beszélt nyelveknél, mint amilyen a magyar is, nagy hátrányt jelent bármilyen szövegértésre, hangfeldolgozásra épülő alkalmazás lokalizálásánál. Annak érdekében, hogy a nemzetközi vívmányok kisebb nyelvi piacokon, így Magyarországon is alkalmazhatóak legyenek, szükséges a mesterséges intelligencia alapú nyelvi megoldások alkalmazása.

Jelen kutatás egy kellően nagy HPC kapacitással (56 petaflops~ 1 millió laptop) rendelkező szolgáltatás segítségével lendületet szeretne ennek adni a hazai fejlesztésnek, és a hosszútávon a lehető legnagyobb hasznos méretű általános nyelvi modell előállítása a cél.

OTP-s NLP „Szuperszámítógép” projekt céljai

Projekt tagjai-Működési modell

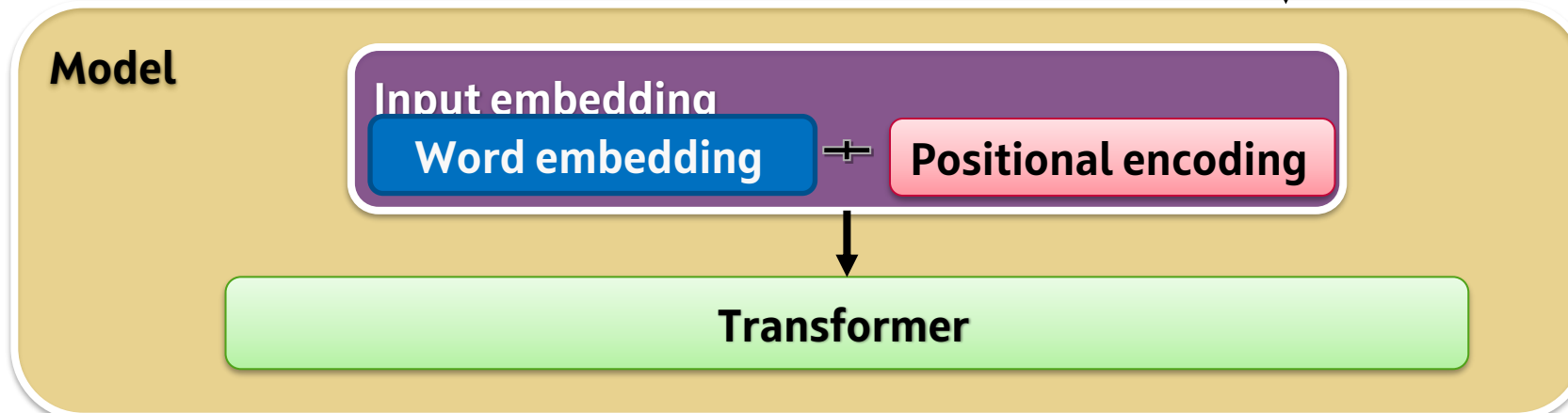
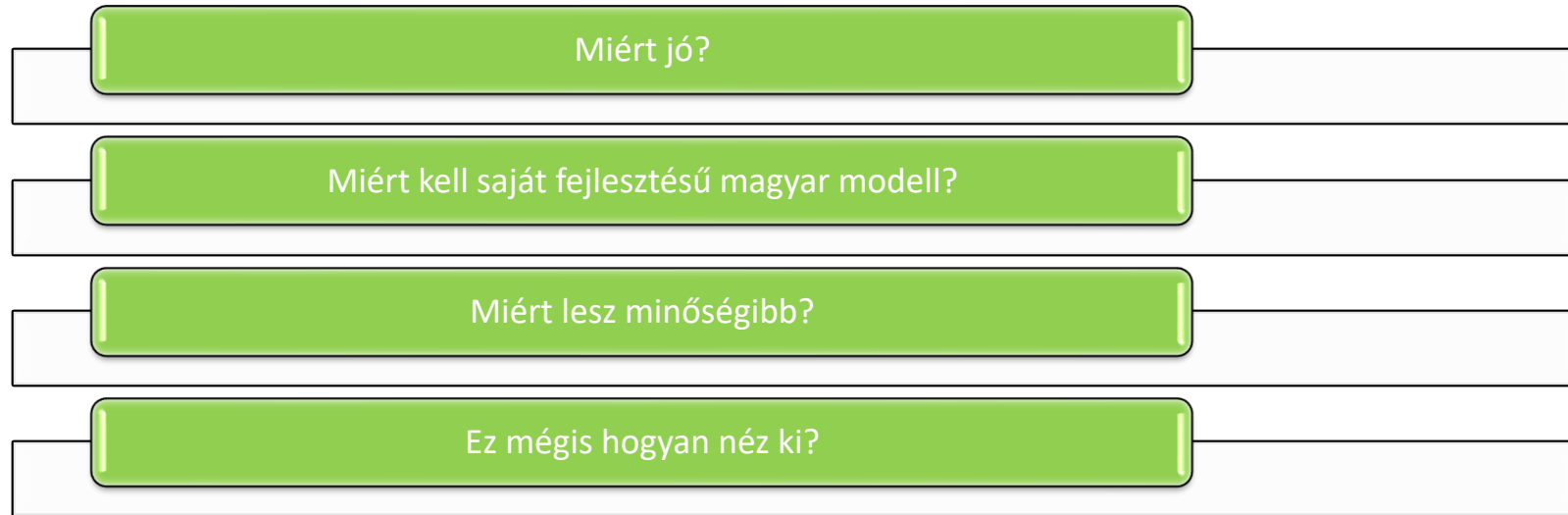
- Pályázat paraméterei
- Nemzeti közkincs egyértelmű cél
- Képesség, tudás építés
- OTP Csoport és rengeteg szereplő segít



A projekt munkacsoport jellegű működéssel indul el, előre definiált, ám rugalmasan alakítható feladatokkal. A tevékenységek előre haladtával javasoljuk a működési modell rendszeres felülvizsgálatát, az igényekhez igazítását és a szükséges kompetencia mix finomhangolását.

ITM és SambaNova-n kívüli fontos külső kapcsolódások jelennek meg: az MNB-vel, a banki MI szabályozási környezet kialakítására és a magyar nyelvtudományi közösség, egyetemeken keresztül, a nyelvészeti kutatások felhasználására és a magyar nyelvi modell kialakítására, K+F tevékenység sikeres elvégzésére.

Nyelvi modellek: miért szükséges a Magyar NLP modell?

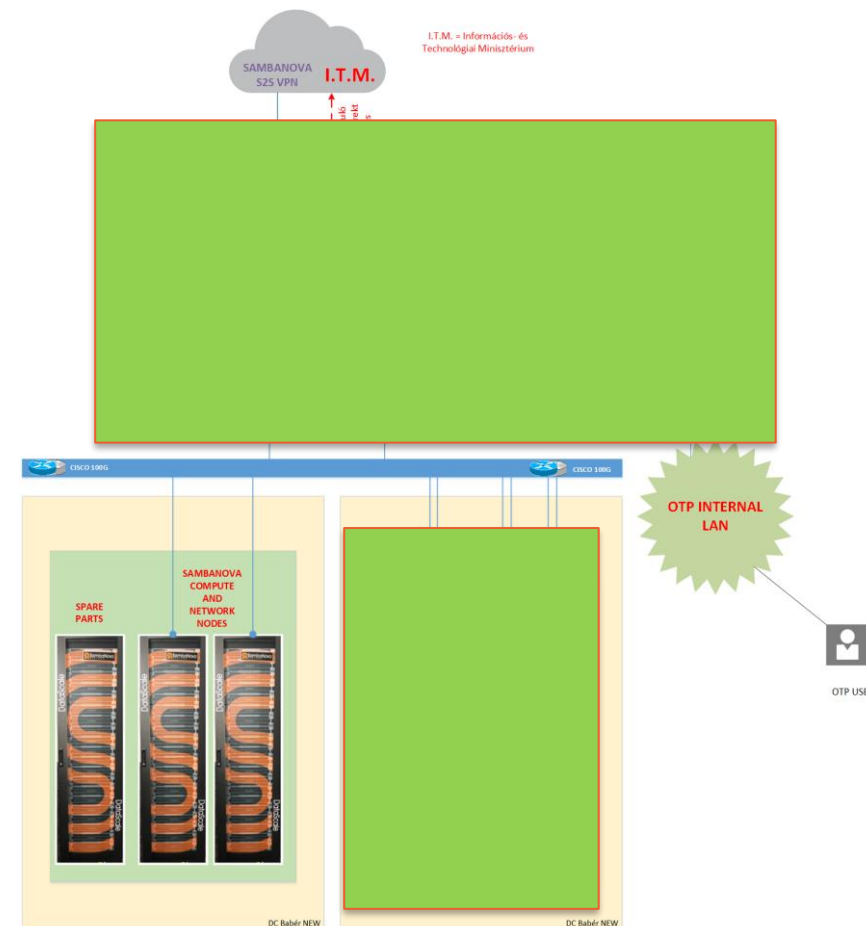


OTP szuperszámítógép választásának szempontjai

OTP szuperszámítógép építés lépései

- Nem általános célú HW-szolgáltatás
- Gyors implementáció volt a szempont
- Covid időben jött és volt HW
- Lehetett egyből dolgozni vele
- CoE Stanford támogatással

- HW szállítás
- 3 rack, 750 kg 1 DB
- Kevlár szőnyeg
- 15-21 KW/h felvétel- mégis hatékony



One SN10-8 system has 8 RDUs. Each RDU has > 300 BF16 TFLOPs. That means that one system has > 2,400 BF16 TFLOPs.

A full Dataflow-as-a-Service rack has 4 x SN10-8 systems. That means that a full rack has > 19,200 BF16 TFLOPs (19.2 PFLOPs)

3 racks has > 57,600 BF16 TFLOPs (57.6 PFLOPs)

Ha banki rendszer lenne:



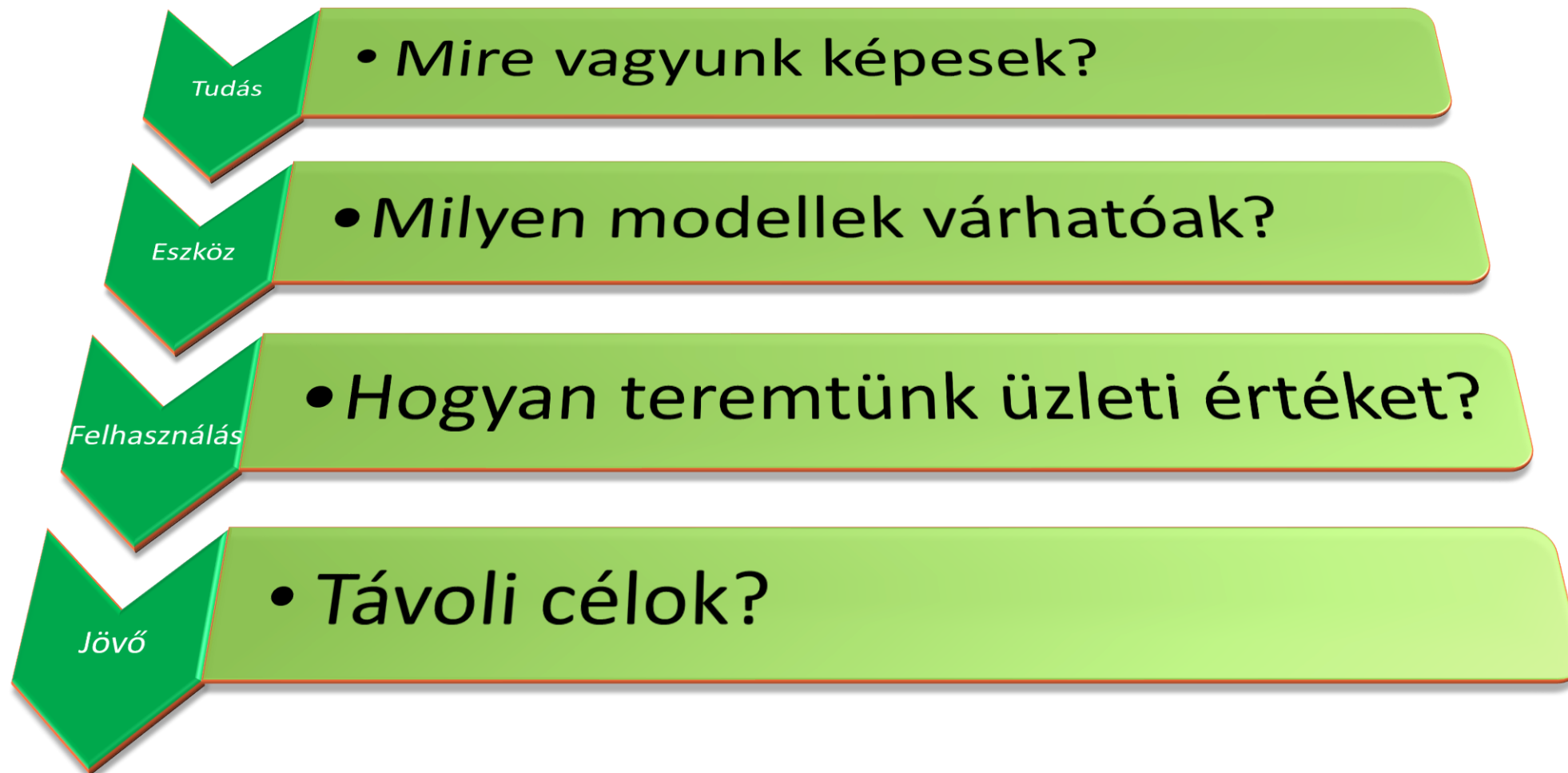
- Alkalmazás ID
- Rto, Rpo, SLA, BIA, Ibtv, MNB fogyasztóvédelem,
- HLD, LLD
- Adatleltár, Adatok BSR , BRD osztály + BRD
- SIEM + SOC, IT monitoring, egyéb pl. speciális elemző eszközök
- BCP, DRP, MAVT
- Szabályzatok, auditok, felügyeleti ügyek, csoport ügyek
-

Nem banki rendszer:



- Terület, fizikai biztonság
- Energia, klíma, + felügyelet
- Adatkapcsolat
- Önálló védelem és IT monitoring és üzem
- Önálló AAA
- Adathordozó ? ekkor sem...





Függelék

The background features several abstract, overlapping geometric shapes in various shades of green. On the left, there are three concentric, semi-circular arcs. To the right, there is a solid circular shape. The overall composition is minimalist and modern.

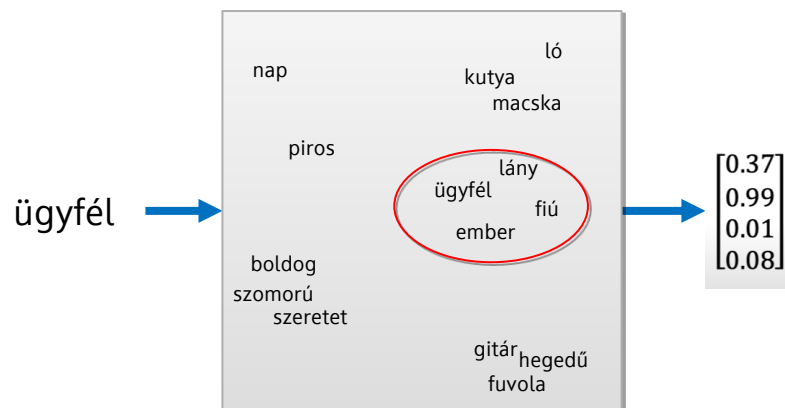
Ez mégis hogyan néz ki?

Input embedding

Word embedding



Positional encoding



Péter nem figyelt, és **elesett**. → Pozíció 5

A király **elesett** a csatában. → Pozíció 3

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}})$$
$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

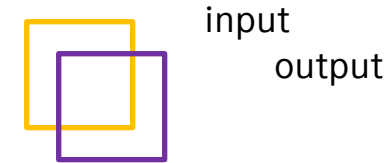


$\begin{bmatrix} 0.42 \\ 0.84 \\ 0.12 \\ 0.81 \end{bmatrix}$

→ Az ügyfél embedding-je, kontextussal

Ez mégis hogyan néz ki?

LM head



	előtérben ügyfeleink pénzügyi biztonsága a												
előtérben	<input type="checkbox"/>	0.02	0.53	0.08	0.14	0.23	ügyfeleink	<input type="checkbox"/>	0.00	1.00	0.00	0.00	0.00
ügyfeleink	<input type="checkbox"/>	0.07	0.00	0.44	0.35	0.14	pénzügyi	<input type="checkbox"/>	0.00	0.00	1.00	0.00	0.00
pénzügyi	<input type="checkbox"/>	0.09	0.21	0.03	0.51	0.16	biztonsága	<input type="checkbox"/>	0.00	0.00	0.00	1.00	0.00
biztonsága	<input type="checkbox"/>	0.14	0.25	0.11	0.01	0.49	a	<input type="checkbox"/>	0.00	0.00	0.00	0.00	1.00

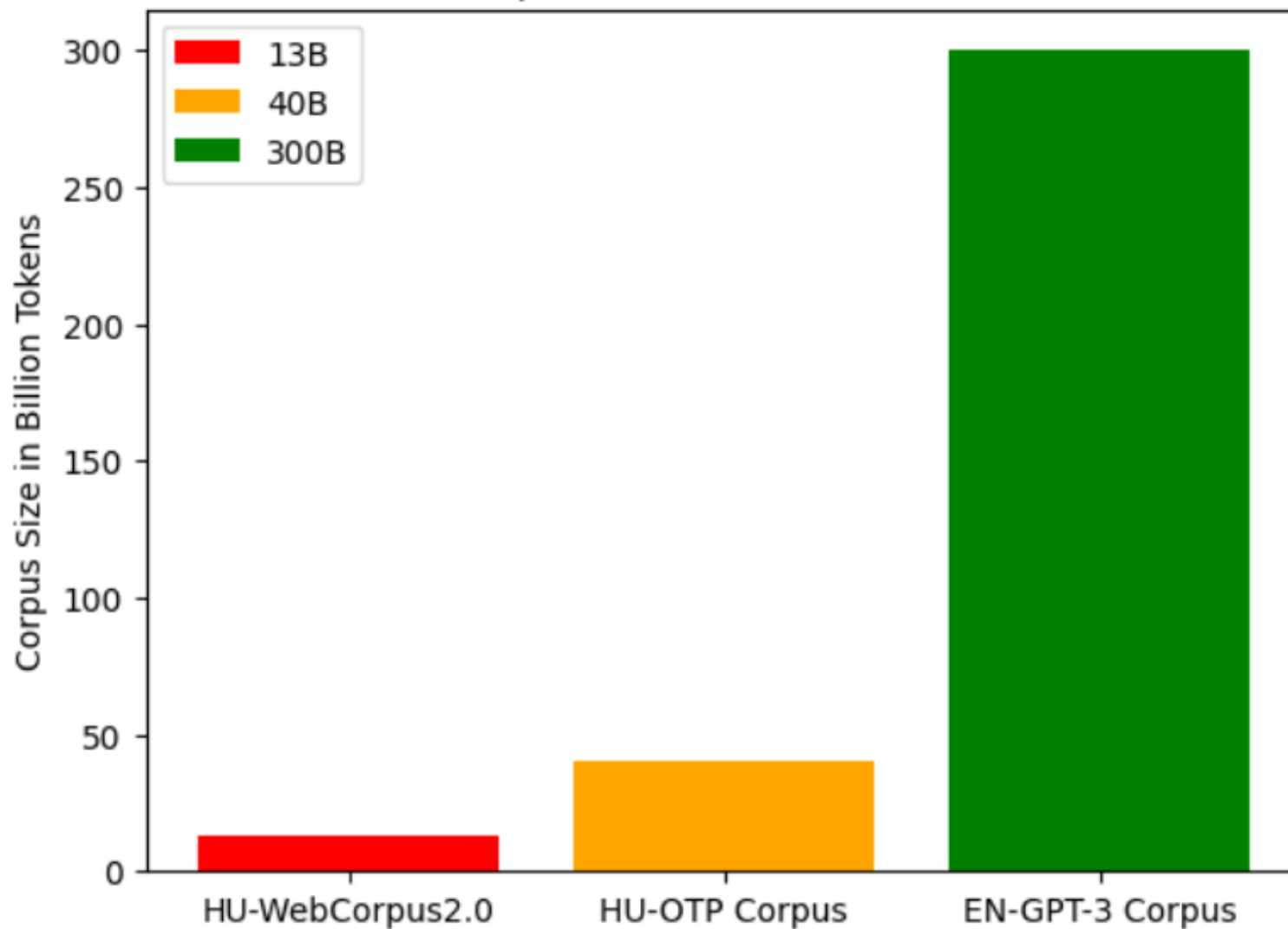
Ez mégis hogyan néz ki?

Classification head



	pozitív	negatív		
előtérben	0.96	0.04		
ügyfeleink	0.65	0.35		
pénzügyi	0.26	0.74		
biztonsága	0.81	0.19		
átlag:	0.67	0.33	1	0

Corpus Size in Billion Tokens



Corpus <> Context Window

1.5B model 1024 token

13B model 2048 token

WebCorpus 2.0 has 13B token:

- 1.5B model **12 397 steps**
- 13B model **6 198 steps**

150k steps would be ideal!

How much data do we need?

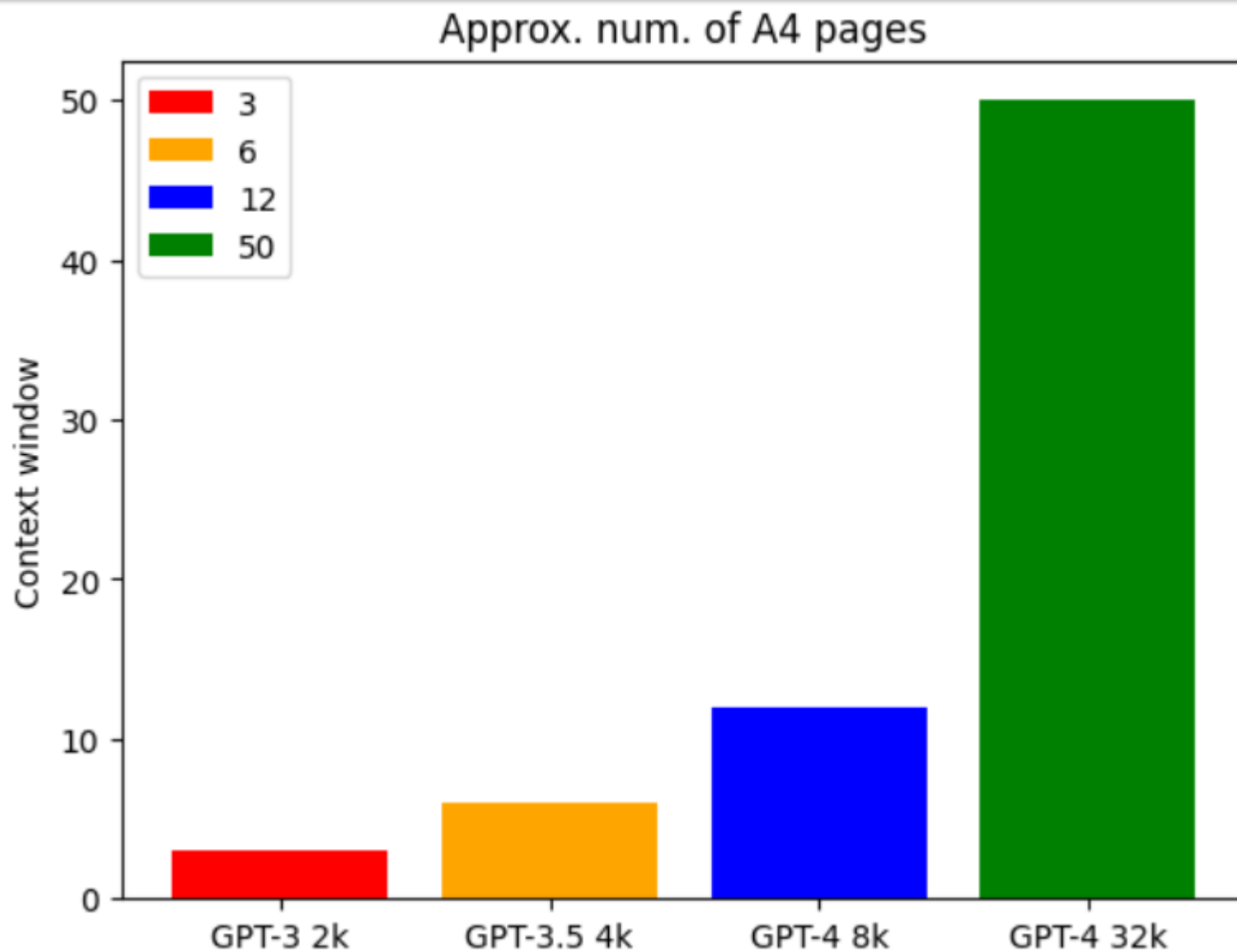
Sambanova's experience:

- 13B model (300k steps) can overperform 175B model (150k steps)

300k steps -> $1024 * 300k * 2048 = 629B$ tokens!!!

150k steps -> **314B tokens!!!**

Context window



Hogyan készülünk a modell tanítására?

„Szótár” létrehozása, előre meghatározott mérettel

Corpus

1	huggingface
1	hugging
1	face
1	hug
1	hugger
2	learning
2	learner
2	learners
1	learn

Képzeljünk el egy szöveges állományt (corpora), amiben a bal oldalon látható **szavak** szerepelnek. **Kék színnel** az előfordulásuk darabszámát jelöljük.

Hogyan készülünk a modell tanítására?

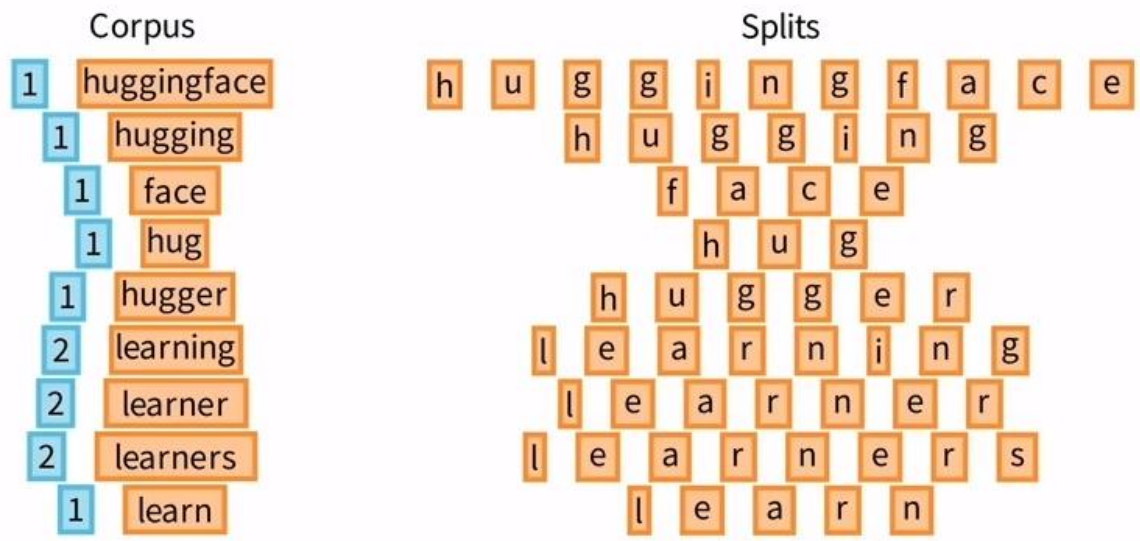
„Szótár” létrehozása, előre meghatározott mérettel



Karakterekre bontjuk (Splits) a szavakat.

Hogyan készülünk a modell tanítására?

„Szótár” létrehozása, előre meghatározott mérettel



Az előforduló karaktereket felvesszük a szótárba.
Van még hely? (megadtuk, mekkora lehet a szótár)
Ha igen, akkor megyünk tovább...

Hogyan készülünk a modell tanítására?

„Szótár” létrehozása, előre meghatározott mérettel



Pairs frequencies

h	+	u	:	4
u	+	g	:	4
g	+	g	:	3
g	+	i	:	2
i	+	n	:	2
n	+	g	:	2
g	+	f	:	1
f	+	a	:	2
a	+	c	:	2
c	+	e	:	2
g	+	e	:	1
e	+	r	:	1
l	+	e	:	1
e	+	a	:	1
a	+	r	:	1
r	+	n	:	1
n	+	i	:	1



A szimpla karakterek után,
hozzadjuk a párokat
szótárhoz.

Hogyan készülünk a modell tanítására?

„Szótár” létrehozása, előre meghatározott mérettel

Corpus

1	huggingface
1	hugging
1	face
1	hug
1	hugger
2	learning
2	learner
2	learners
1	learn

Splits

h	u	g	g	i	n	g	f	a	c	e
h	u	g	g	i	n	g				
		f	a	c	e					
		h	u	g						
	h	u	g	g	e	r				
l	e	a	r	n	i	n	g			
l	e	a	r	n	e	r				
l	e	a	r	n	e	r	s			
		l	e	a	r	n				

Pairs frequencies

h	+	u	:	4
u	+	g	:	4
g	+	g	:	3
g	+	i	:	2
i	+	n	:	3
n	+	g	:	3
g	+	f	:	1
f	+	a	:	2
a	+	c	:	2
c	+	e	:	2
g	+	e	:	1
e	+	r	:	3
l	+	e	:	4
e	+	a	:	4
a	+	r	:	4
r	+	n	:	4
n	+	i	:	1
n	+	e	:	2
r	+	s	:	1

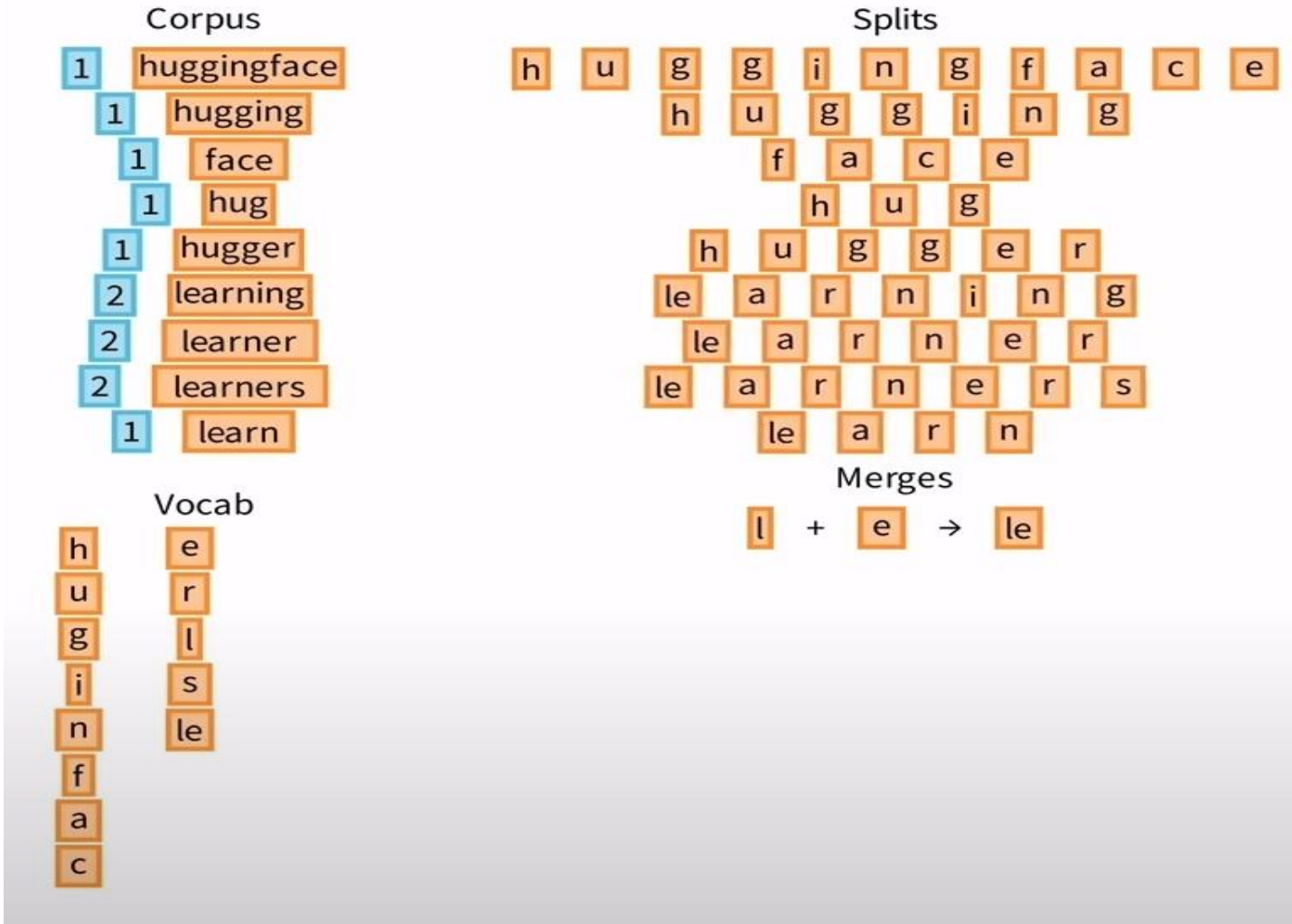
Vocab

h	e
u	r
g	l
i	s
n	
f	
a	
c	

Amelyik párból a legtöbb van, cseréljük le a két szimpla karakter a párra -> a fenti Splits listában.

Hogyan készülünk a modell tanítására?

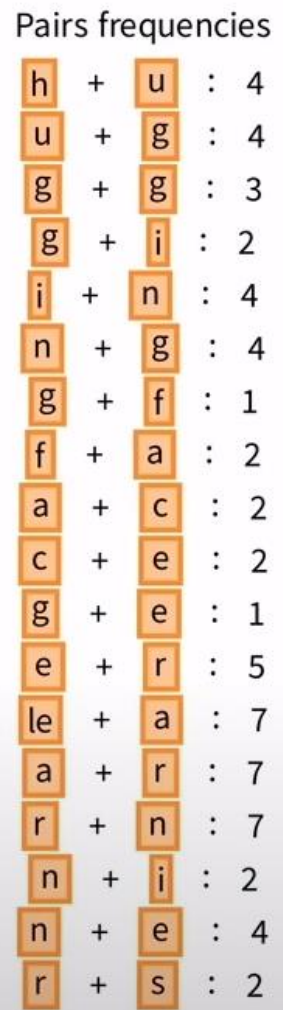
„Szótár” létrehozása, előre meghatározott mérettel



l e -> le

Hogyan készülünk a modell tanítására?

„Szótár” létrehozása, előre meghatározott mérettel



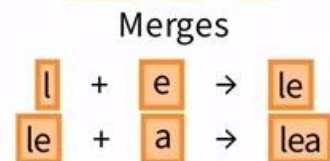
Keressünk újra token párokat és számoljuk meg az előfordulásuk gyakoriságát

Hogyan készülünk a modell tanítására?

„Szótár” létrehozása, előre meghatározott mérettel



Pairs frequencies



Hogyan készülünk a modell tanítására?

Corpus

1 huggingface
1 hugging
1 face
1 hug
1 hugger
2 learning
2 learner
2 learners
1 learn

Splits

hug g in g f a c e
hug g in g
f a c e
hug
hug g er
learn in g
learn er
learn er s
learn

Pairs frequencies

hug + g : 3
g + in : 2
in + g : 4
g + f : 1
f + a : 2
a + c : 2
c + e : 2
g + er : 1
learn + in : 2
learn + er : 4
er + s : 2

Vocab

h e
u r er
g l hu
i s hug
n le in
f lea
a lear
c learn

Merges

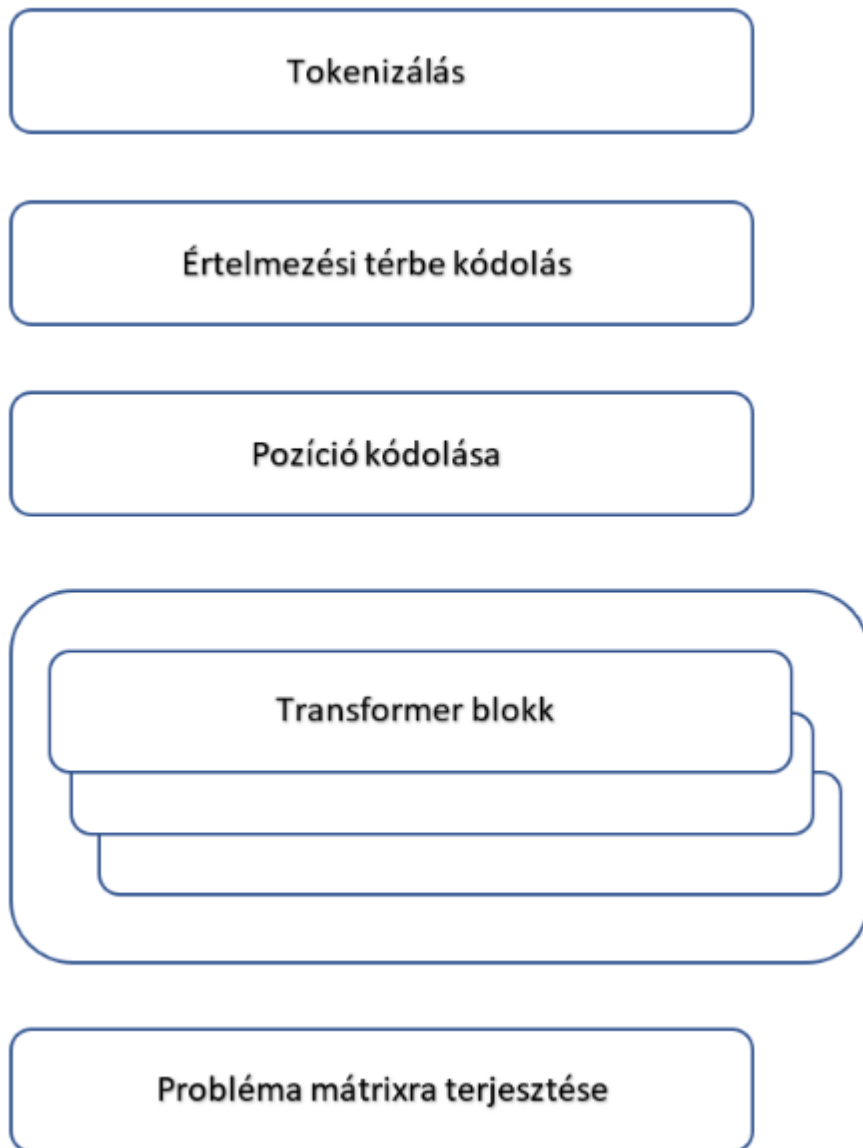
l + e → le
le + a → lea
lea + r → lear
lear + n → learn
e + r → er
h + u → hu
hu + g → hug
i + n → in

Addig csináljuk, amíg el nem fogy a hely a szótárban.

Mekkora legyen a magyar szótár?

Megtaláltuk, ez az egyik kutatási eredményünk

Hogyan készülünk a modell tanítására?



Következőnek
átnézhetjük az
értelmezési térbe
kódolást.

Köszönjük a figyelmet!

The background features several overlapping, semi-transparent green geometric shapes. On the left, there are two concentric, thick curved lines that resemble a stylized 'C' or a partial circle. To the right of these, there is a solid green circle. The overall composition is minimalist and modern, using a monochromatic green color palette.

Kérdések?

The background features several overlapping, semi-transparent green geometric shapes. On the left, there are three concentric, semi-circular arcs of varying radii. To the right of these arcs is a solid green circle. The overall composition is minimalist and modern, set against a solid green background.