



Seznam LLM v produkcí



Marek Šimůnek

Tech lead LLM týmu



Agenda

- Proč děláme vlastní Seznam LLM => SeLLMa => Šelma?
- Škálování inferencí
- Garance provozního trafficu
- Offline (Batchová) cesta



O mně



40,000,000,000 klíčů v databázi
300,000,000 stažení denně

- Big data
- MLOps
- LLM

1000 fyzických serverů
15 PB úložného prostoru
65 TB paměti

Platforma pro 5 výzkumnických týmů
25+ lidí
100+ modelů

Seznam LLM co umí česky



Lupa.cz » Seznam.cz chystá vlastní umělou inteligenci. V češtině už je o něco lepší než GPT-3.5

Seznam.cz chystá vlastní umělou inteligenci. V češtině už je o něco lepší než GPT-3.5

JAN SEDLÁK | 17. 1. 2024 | Doba čtení: 4 minuty

9 NOVÝCH NÁZ

Lupa.cz » Jazykový model Seznamu se jmenuje Šelma, firma pracuje na zapojení AI do služeb

Jazykový model Seznamu se jmenuje Šelma, firma pracuje na zapojení AI do

Zprávy » Byznys » Byznys | Rozhovory » Vyhledávání na Seznamu se dramaticky změní, říká Pav...

Vyhledávání na Seznamu se dramaticky změní, říká Pavel Zima z vedení firmy

UMĚLÁ INTELIIGENCE – 07. 3. 2024 – 4 min čtení

Nejrychlejších sedm měsíců života, říká žena, která vede vývoj generativní umělé inteligence Seznamu

Seznam.cz si jako česká internetová jednička nechce nechat příležitost v umělé inteligenci ujít. Investuje do ní desítky milionů a učí ji česky.



LLM v Seznamu



**Výborná čeština
modelu**



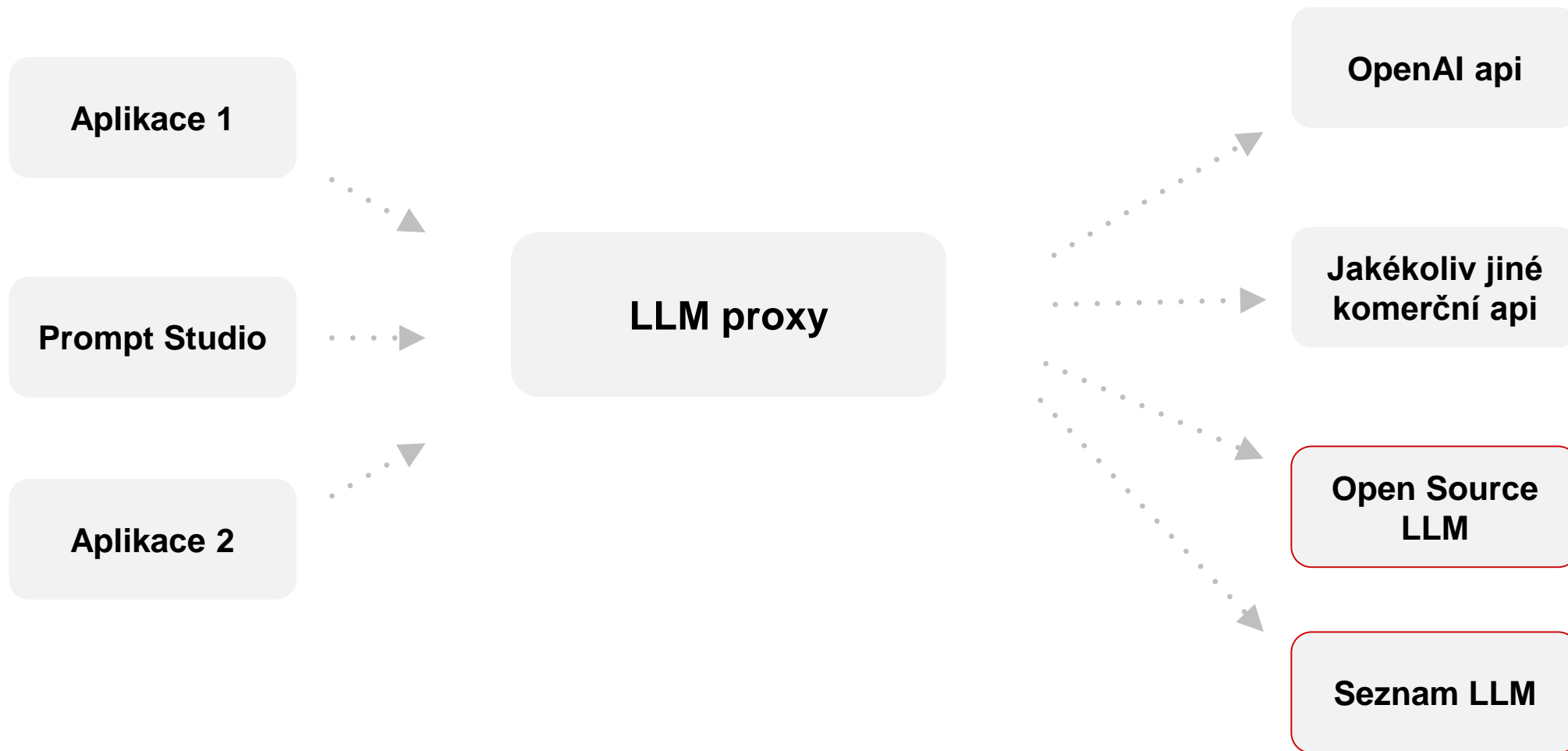
Provoz v naší režii



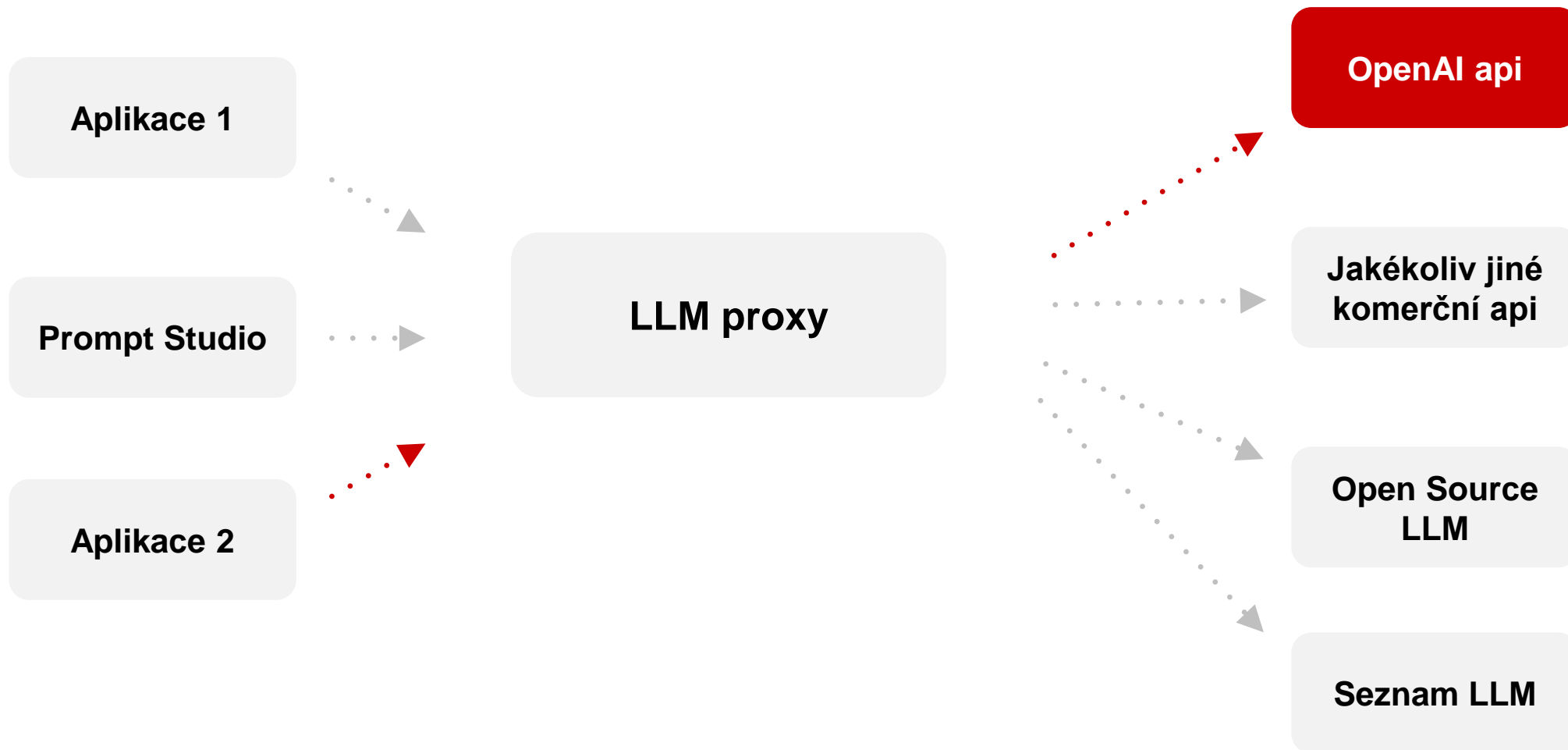
**Bezpečnost pro
uživatele**



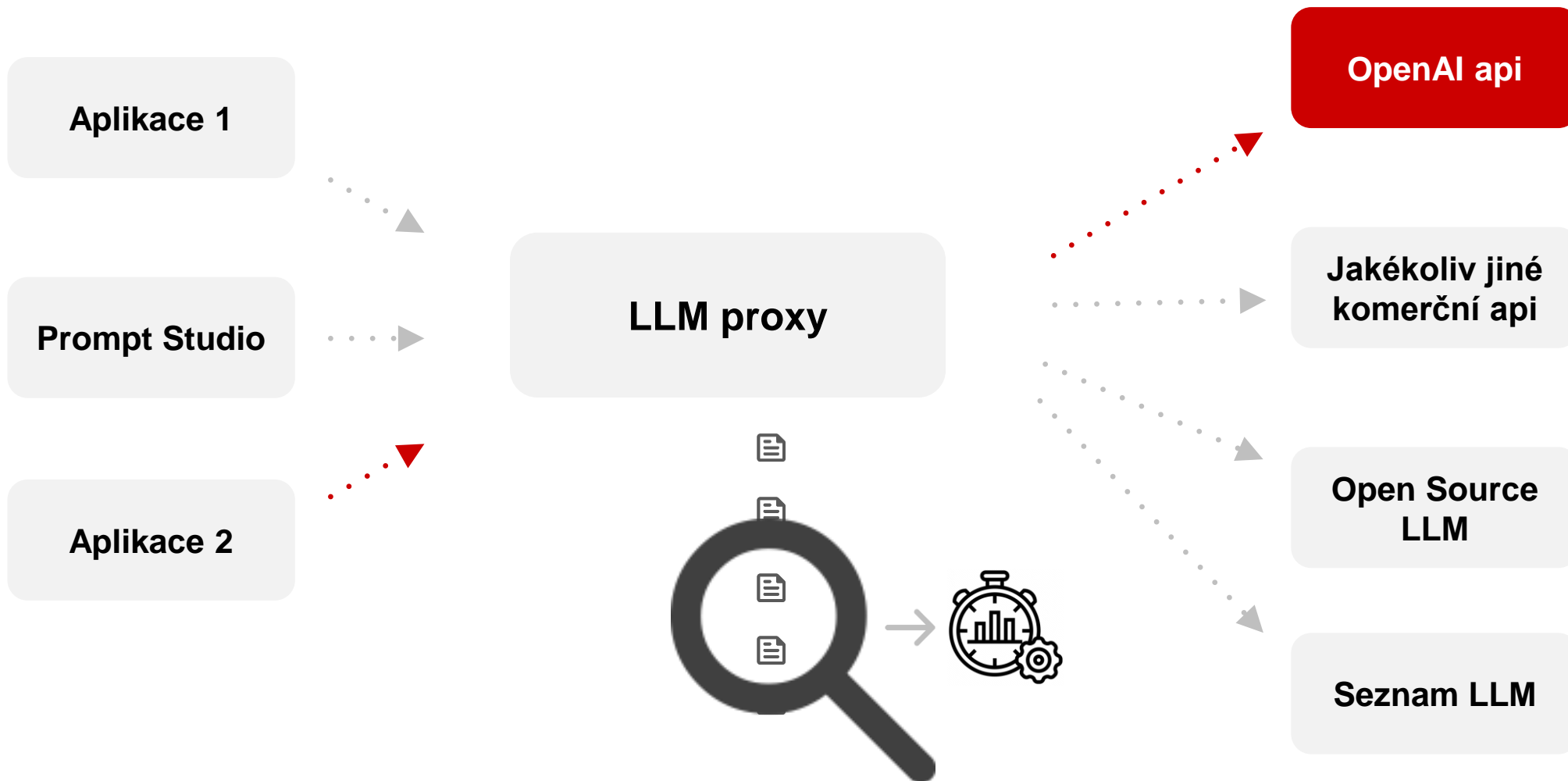
LLM proxy



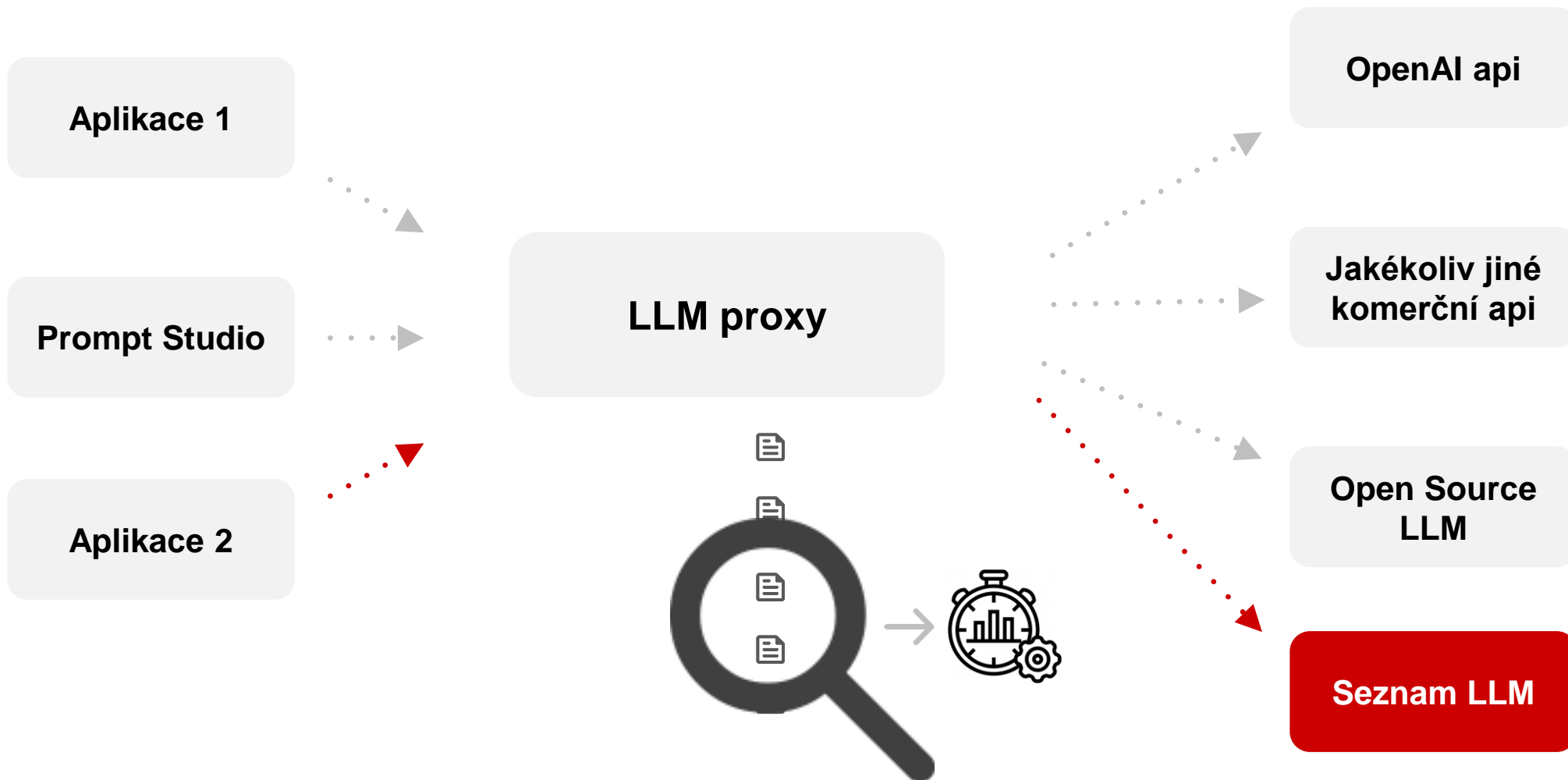
LLM proxy



LLM proxy

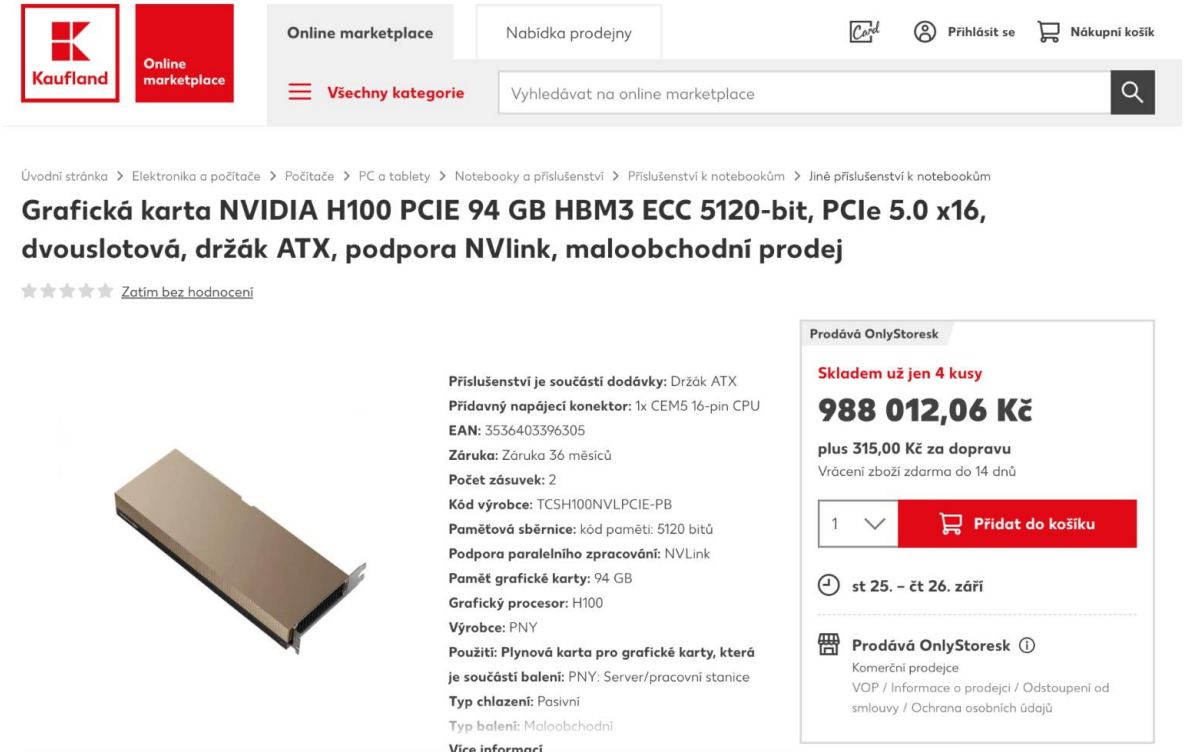


LLM proxy



Požadavky na produkční řešení

- 1 GPU - cena cca 800k - 1M Kč
- proměnlivý traffic během dne
- 10 až 20 modelů v provozu (aktuálně 5)
- garance trafficu pro více uživatelů/aplikací
- umět dodržovat určité latence

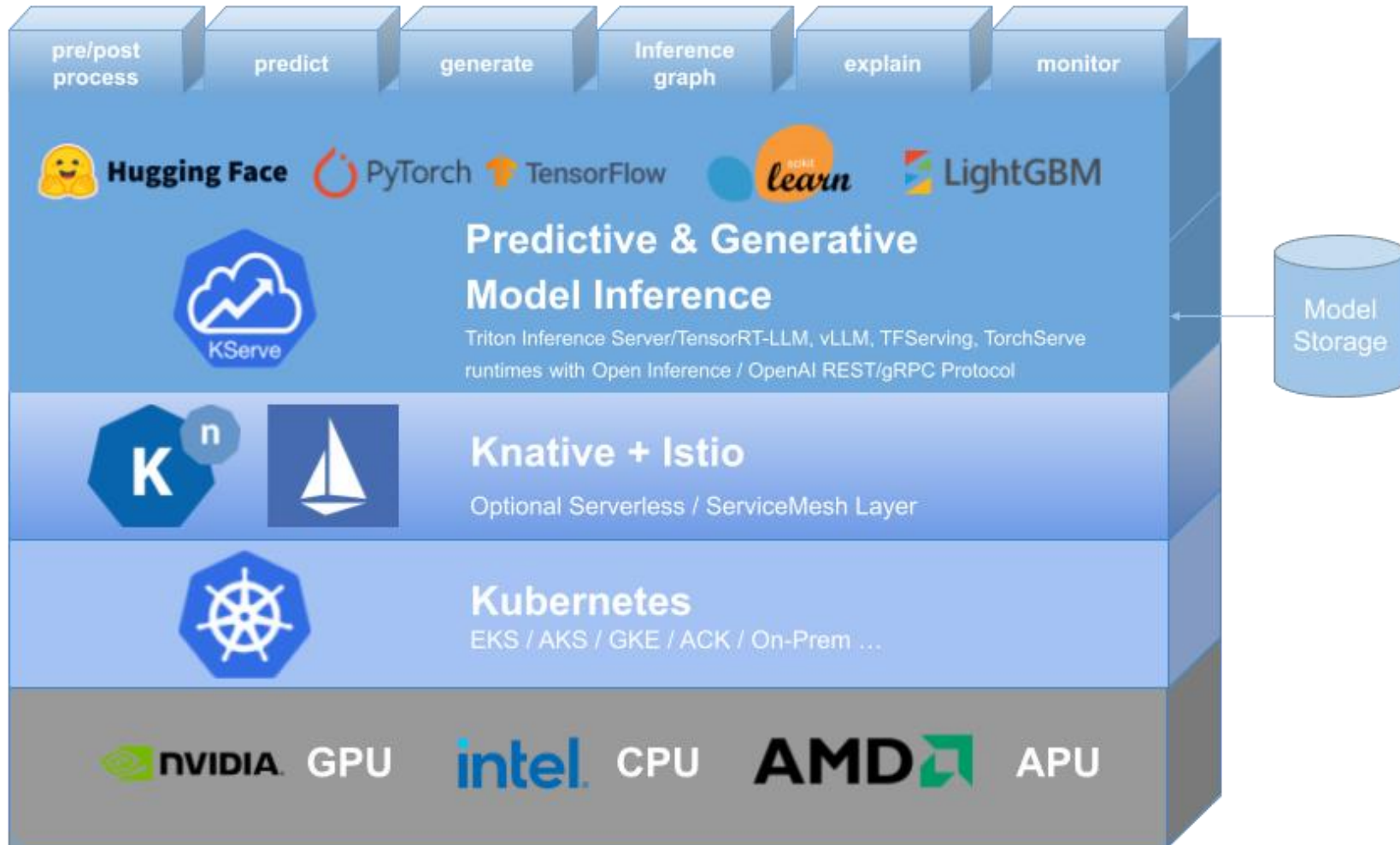


The screenshot shows the product page for an NVIDIA H100 GPU on the Kaufland online marketplace. The page includes the following information:

- Product Title:** Grafická karta NVIDIA H100 PCIE 94 GB HBM3 ECC 5120-bit, PCIe 5.0 x16, dvouslotová, držák ATX, podpora NVLink, maloobchodní prodej
- Price:** 988 012,06 Kč (plus 315,00 Kč za dopravu)
- Availability:** Skladem už jen 4 kusy
- Shipping:** Vrácení zboží zdarma do 14 dnů
- Product Details:**
 - Příslušenství je součástí dodávky: Držák ATX
 - Přídavný napájecí konektor: 1x CEM5 16-pin CPU
 - EAN: 3536403396305
 - Záruka: Záruka 36 měsíců
 - Počet zásuvek: 2
 - Kód výrobce: TCSH100NVLPcie-PB
 - Paměťová sběrnice: kód paměti: 5120 bitů
 - Podpora paralelního zpracování: NVLink
 - Paměť grafické karty: 94 GB
 - Grafický procesor: H100
 - Výrobce: PNY
 - Použití: Plynová karta pro grafické karty, která je součástí balení: PNY: Server/pracovní stanice
 - Typ chlazení: Pasivní
 - Typ balení: Maloobchodní
 - Více informací



KServe



KServe: jak řešíme škálování

```
apiVersion: serving.kserve.io/v1beta1
```

```
kind: InferenceService
```

```
metadata:
```

```
  name: torch-service
```

```
spec:
```

```
  predictor:
```

```
    scaleTarget: 10
```

```
    scaleMetric: concurrency
```

```
  model:
```

```
    modelFormat:
```

```
      name: pytorch
```

```
    storageUri: pvc://some-persistent-volume-claim
```

metrika podle, které škálovat



jaký je formát modelu



odkud model načíst

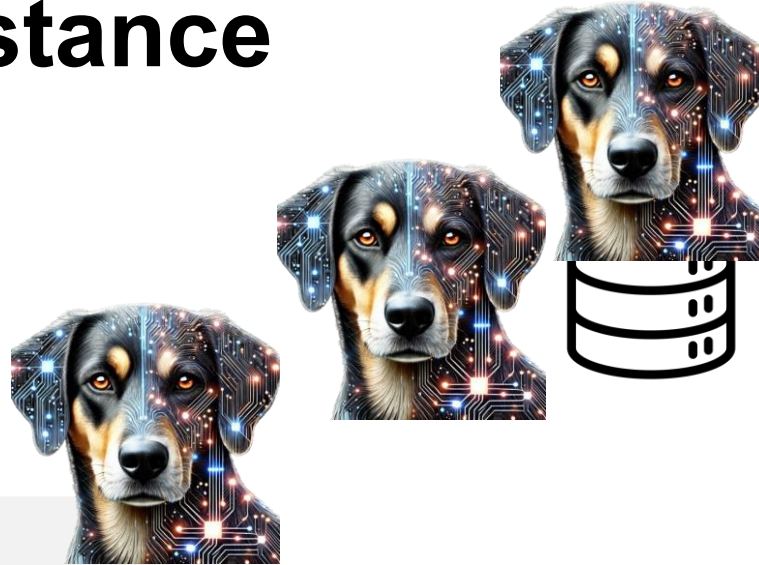


KServe: jak řešíme škálování

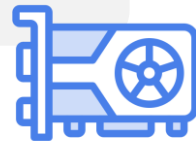
```
apiVersion: serving.kserve.io/v1beta1
kind: InferenceService
metadata:
  name: torch-service
spec:
  predictor:
    maxReplicas: 5
    minReplicas: 0
    scaleTarget: 10
    scaleMetric: concurrency
  model:
    modelFormat:
      name: pytorch
    storageUri: pvc://network-storage-gate-pvc
```

Nastartování nové instance

140GB



Inference service 1



Doba naběhnutí:

- sellma-70B: 4-5min
- sellma-7B: 1-2min

Applications use cases



Online

Chci odpověď co nejdřív



Offline

Mám hodně requestů, ale odpověď
počká



Applications use cases



Online

Balancování mezi dobou odpovědí a celkovou propustností systému.



Offline

Nejjednodušší: Zaměření na maximální propustnost



Applications use cases



Online

Balancování mezi dobou odpovědi a celkovou propustností systému.



Offline

Nejjednodušší: Zaměření na maximální propustnost

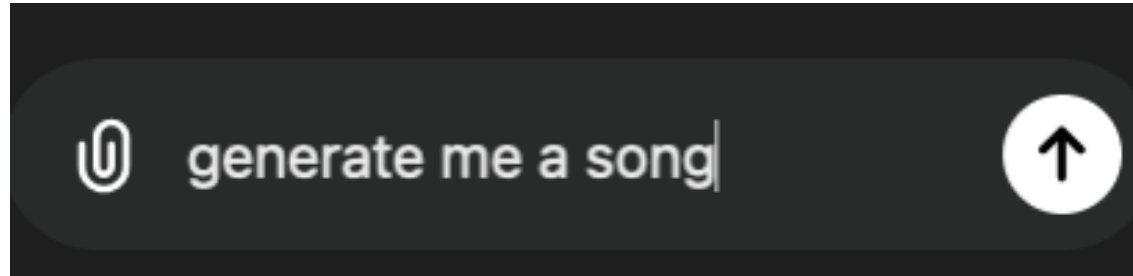
Streaming

Sequential



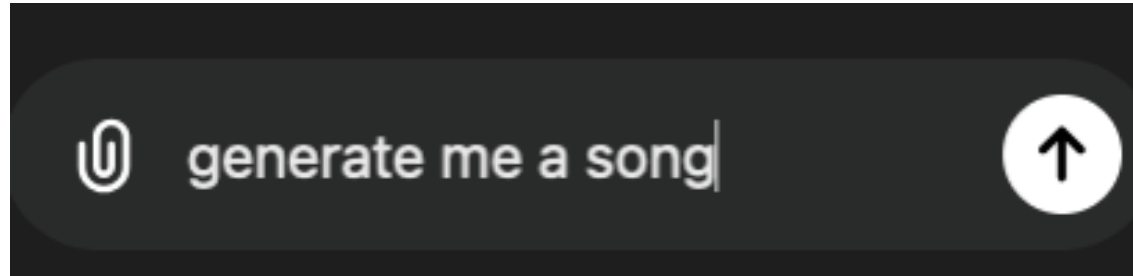
Vysvětlení inference pro streaming

Příklad:



Vysvětlení inference pro streaming

Příklad:



Prefill: 1.14s, 5 input tokens, 1 output token

generate me a song

Decoding: 1.62s, 33 output tokens

Could you please provide me with some specifics

Prefill, first-response	Decoding							
User: generate me a song. AI: <i>Could</i>	you	please	me	with	some	specifics	for	your

prompt

user prompt

Vysvětlení inference pro streaming

Prefill = time to first token (TTFT)

- Načtení celého promptu
- Čas od načtení promptu k vrácení prvního tokenu
- Závislé na počtu vstupních tokenů
- Limitovaný rychlostí výpočtu na GPUčku

Decoding = inter-token latency (ITL)

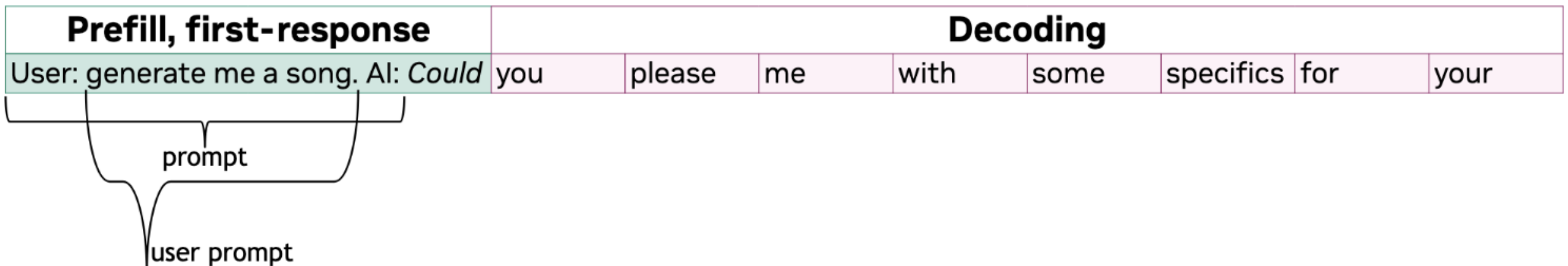
- Generování výstupu token po tokenu
- ITL závisí jak na velikosti vstupu a tak i velikosti výstupu
- Limitovaný rychlostí GPU paměti

Prefill: 1.14s, 5 input tokens, 1 output token

generate me a song

Decoding: 1.62s, 33 output tokens

Could you please provide me with some specifics



Applications use cases



Online

Balancování mezi dobou odpovědi a celkovou propustností systému.

Streaming

metriky:

- Time to first token
- Inter token latency

Sequential

metriky:

- End to end latency (request latency)



Offline (Batch)

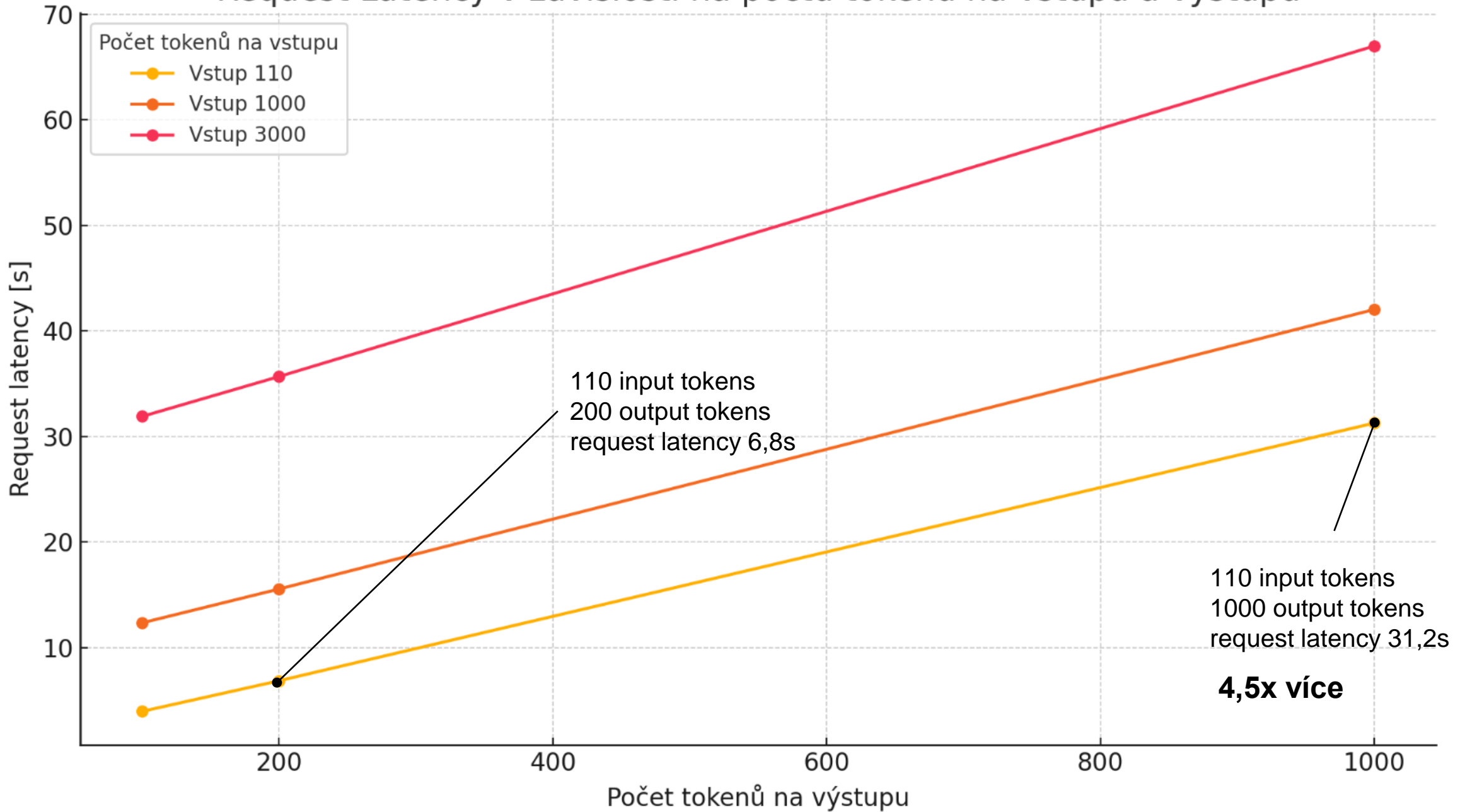
Nejjednodušší: Zaměření na maximální propustnost

metriky:

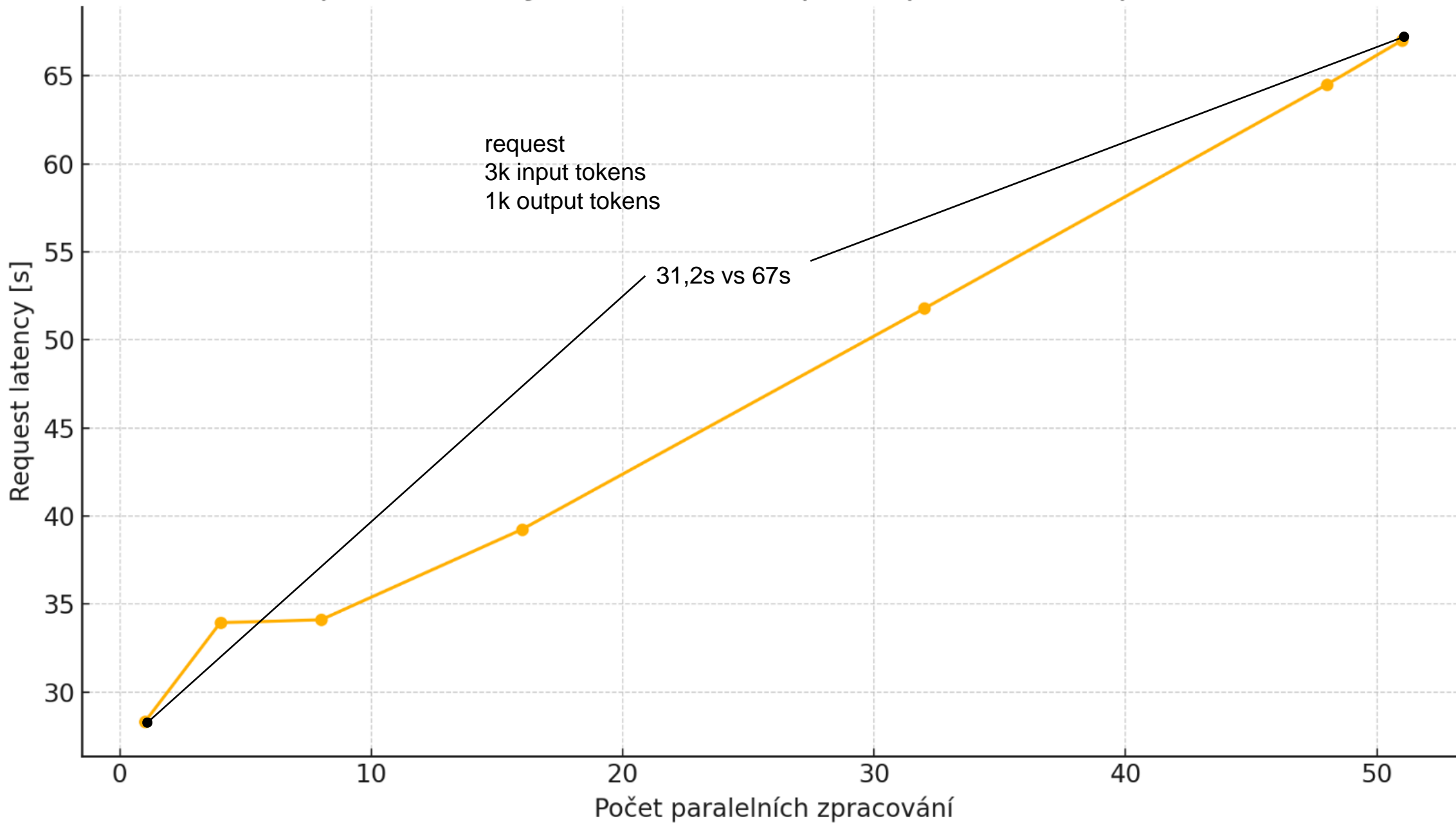
- request_throughput
- output_token_throughput



Request Latency v závislosti na počtu tokenů na vstupu a výstupu



Request Latency v závislosti na počtu paralelních zpracování



Porovnání

Sellma-70B, 1k input tokens, 1k output tokens, H100 PCIe

metric	fp16-4gpu	2x fp8-2gpu	fp8-4gpu
request_latency [s]	60.9	33.3	63.6
request_throughput [req/s]	2.13	3.06	4.01
time_to_first_token [s]	0.28	0.073	0.269
inter_token_latency [s]	0.05	0.028	0.054



Porovnání

Sellma-70B, 1k input tokens, 1k output tokens, H100 PCIe

metric	fp16-4gpu	2x fp8-2gpu	fp8-4gpu
request_latency [s]	60.9	33.3	63.6
request_throughput [req/s]	2.13	3.06	4.01
time_to_first_token [s]	0.28	0.073	0.269
inter_token_latency [s]	0.05	0.028	0.054

0 až 0.1s
pocitově
instantní



Porovnání

Sellma-70B, 1k input tokens, 1k output tokens, H100 PCIe

metric	fp16-4gpu	2x fp8-2gpu	fp8-4gpu
request_latency [s]	60.9	33.3	63.6
request_throughput [req/s]	2.13	3.06	4.01
time_to_first_token [s]	0.28	0.073	0.269
inter_token_latency [s]	0.05	0.028	0.054

0.1s rychlost
lidského čtení



Nasazování aplikací

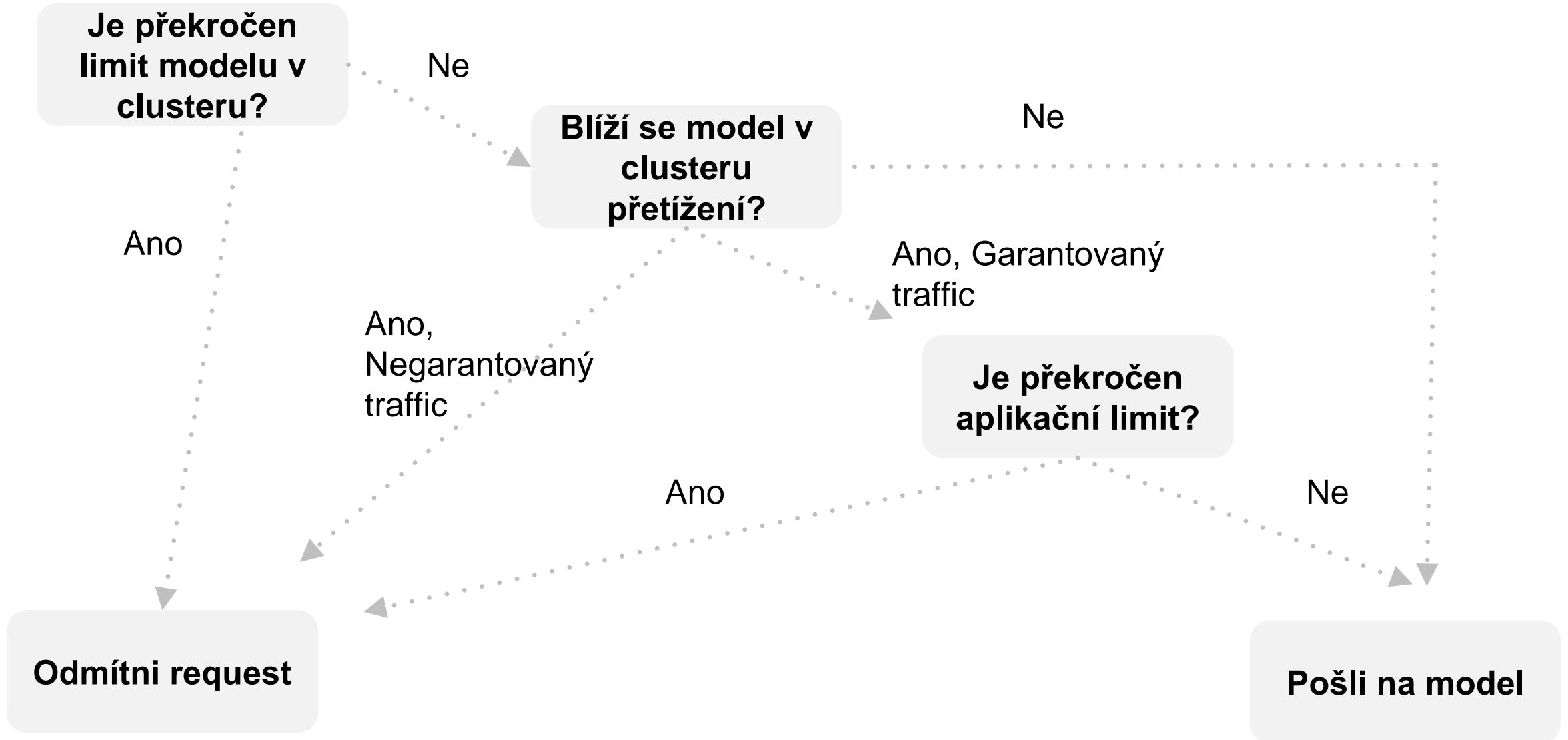
- Jaký model?
- Kolik v průměru tokenů na vstupu a výstupu
 - do velikosti vstupu se počítá i např historie chatu atd.
- Kolik requestů ve špičce
- Jaká je očekávaná doba odpovědi



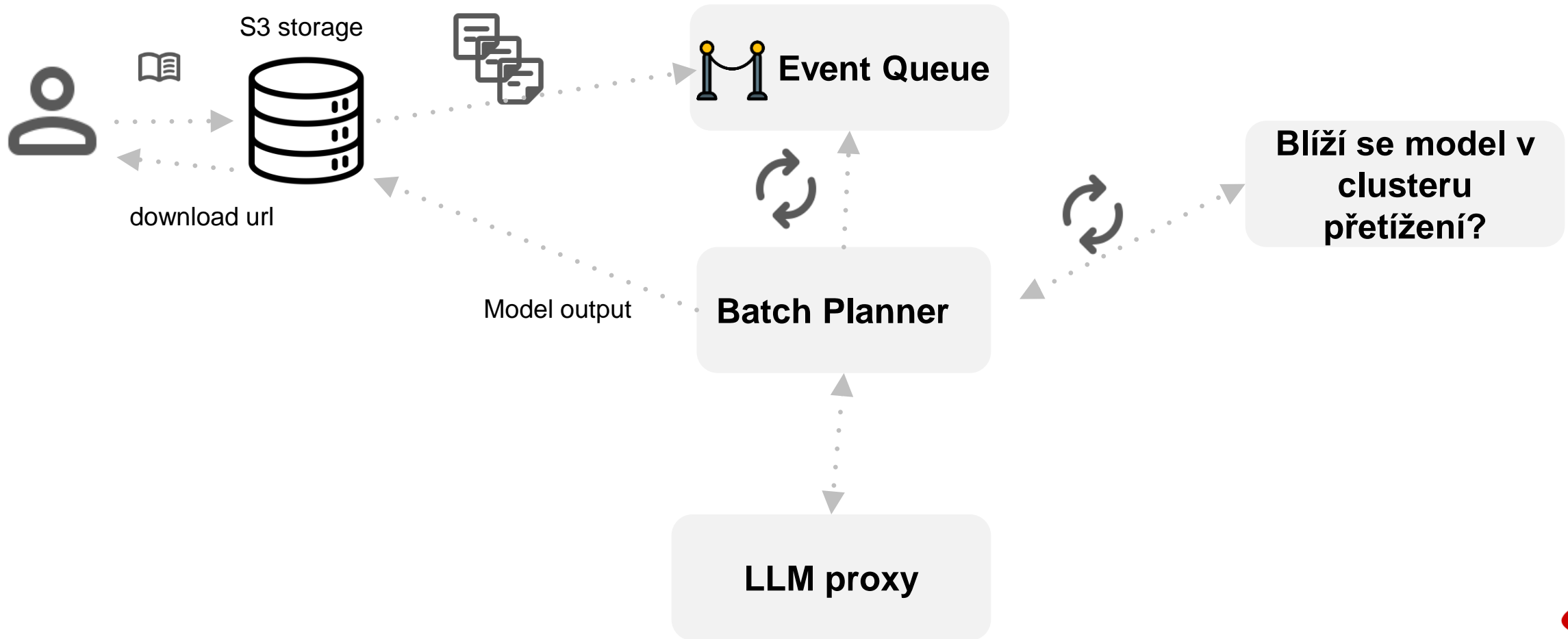
Garance pro provozní aplikace

Garantovaná metrika: počet paralelních requestů (concurrency nebo in-flight req)

Garance pro provozní aplikace



Offline (batch) inference



Prefix caching - sdílený system prompt

Request 1

Udělej souhrn následujícího textu. Myšlenky ****SE NESMÍ**** opakovat. Udělej výstup formou krátkého souvislého odstavce. Odstavec by měl mít maximálně 5 vět. U každé věty si uchovej zdroj, odkud pochází. V každé větě zvýrazni nejdůležitější slova pomocí HTML značky ****. Vrať ****POUZE**** RFC8259 JSON pole v následujícím formátu bez odchylek.

Text:

0 podporu EU žádají dva české těžební projekty...

Request 2

Udělej souhrn následujícího textu. Myšlenky ****SE NESMÍ**** opakovat. Udělej výstup formou krátkého souvislého odstavce. Odstavec by měl mít maximálně 5 vět. U každé věty si uchovej zdroj, odkud pochází. V každé větě zvýrazni nejdůležitější slova pomocí HTML značky ****. Vrať ****POUZE**** RFC8259 JSON pole v následujícím formátu bez odchylek.

Text:

Ve snaze oddálit stárnutí a udržet se v kondici držíme diety...



Speculative decoding

```
def below_threshold(l: list, t: int):  
    """Return True if all numbers in the list l are below threshold t.  
    >>> below_threshold([1, 2, 4, 10], 100)  
    True  
    >>> below_threshold([1, 20, 4, 10], 5)  
    False  
    """  
    if isinstance(l, list):  
        return True  
    else:  
        if t <= l < below_threshold(l, t):  
            return True  
        else:  
            # If the first l element of l is an integer, then it is  
            # the whole range of integers.  
            if not isinstance(l[0], list):  
                return True  
            else:  
                # If the first l element of l is a str, then it is  
                # the whole string.  
                if hasattr(l, 'findlen'):  
                    return findlen(l)  
                return False  
  
def thresh(t: int, max: int) -> int:  
    """Return
```

Figure 3: A visualization of the origin of tokens in an example T=1 HumanEval completion. Green background originates with the N-gram draft² model, blue the draft model, and red the oracle model. (Of course, all tokens are eventually checked by the oracle model.) Obvious tokens – like whitespace – are preferentially accelerated relative to difficult ones.

[\[link to paper\]](#)



Shrnutí

- máme nové modely v llm-proxy
- umíme garantovat aplikacím traffic
- pro větší utilizaci jsme vyrobili i offline (batchovou cestu)
- máme v plánu zrychlovat (kvantizace, prefix caching, speculative decoding atd..)



proč polníček

Internet Obrázky Zboží Mapy Vídea Zprávy Firmy Slovník

- Polníček je **neskutečně výživný** a chutí byste si ho snadno spletli s hlávkovým salátem.¹
- Jelikož ho lze **snadno pěstovat** a má významné nutriční vlastnosti, rozhodně byste ho neměli přehlížet, protože může ozdobit i ochutit vaše pokrmy.¹
- Polníček obsahuje **vitamin C, E, B, provitamin A, kyselinu listovou, hořčík a železo**, což z něj činí cenný zdroj živin.²
- V období vegetačního klidu nám může **poskytovat tolik důležité vitamíny**, což je jeho výjimečnost.²
- Polníček je **jednoduchý důkaz** toho, že i v zimě si můžeme doma nachystat čerstvý a velice chutný salát.²

Polníček: Jak jej pěstovat a čím je...
kupi.cz 1

Polníček – Sezónní potraviny
sezonka.cz 2

K Polníček: Jak jej pěstovat
kupi.cz

Před 14 dny · Vě
polní? Kde se vze
článku!



Pastrňák svatba

Internet Obrázky Zboží Mapy Vídea Zprávy Firmy Slovník

Zprávy na téma: Pastrňák svatba

Český hokejový útočník David Pastrňák se oženil se švédskou partnerkou Rebecou Rohlssonovou na chorvatském ostrově Hvar, kde vychovávají roční dceru Freyu Ivy.^{1, 2, 3} Pastrňák se na Instagramu pochlubil fotkou ze svatby, která byla třídenním luxusním veselím, na němž nešetřil a které stálo miliony.^{4, 5}

Tento text je vygenerovaným shrnutím tématu

Adele

Adele – Wikipedie
cs.wikipedia.org/wiki/adele

Adele Laurie Blue Adkins, MBE (* 5. května 1988 Londýn, Spojené království) je anglická zpěvačka a skladatelka známá jen jako Adele.

Kariéra Soukromý život Diskografie Turné Odkazy

Mohlo by vás zajímat

Jaké je celé jméno zpěvačky Adele?

Jaké je debl

Proč se hledá Nehoda na Rakovnicku?

ošlo u Lubné na Rakovnicku k tragické nehodě, při které řidička sjela ze silnice a převrátila auto na střechu.^{1, 2, 3} Řidička na místě o dvě děti, které s ní cestovaly, byly zraněny a převezeny do
^{1, 4, 5}

2. idnes.cz 3. extra.cz 4. iprima.cz 5. blesk.cz

erovaným shrnutím tématu



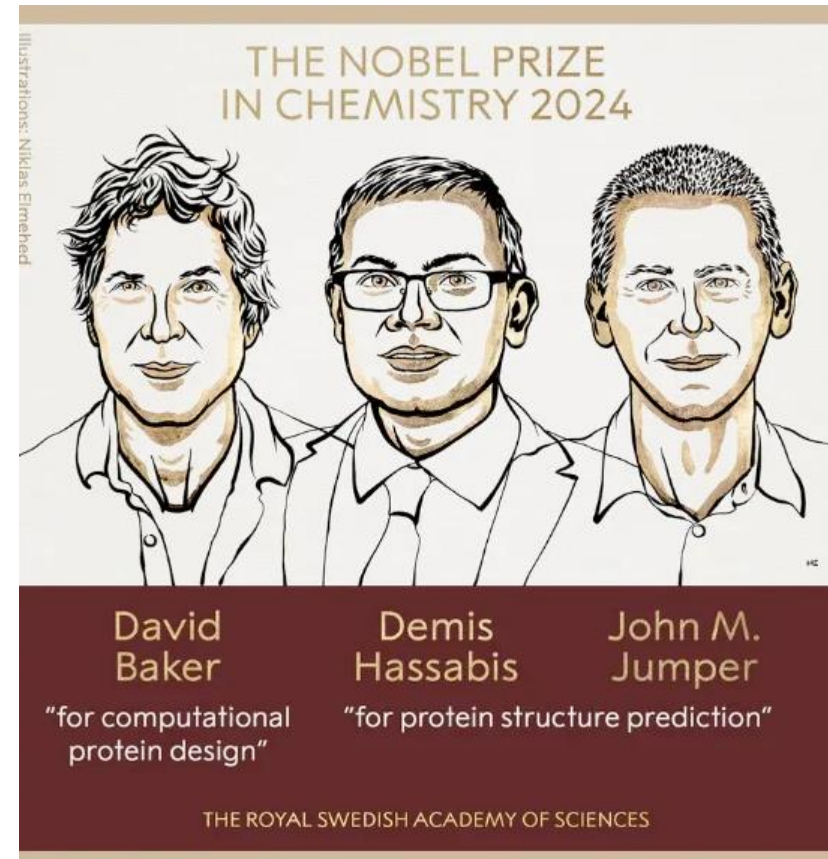
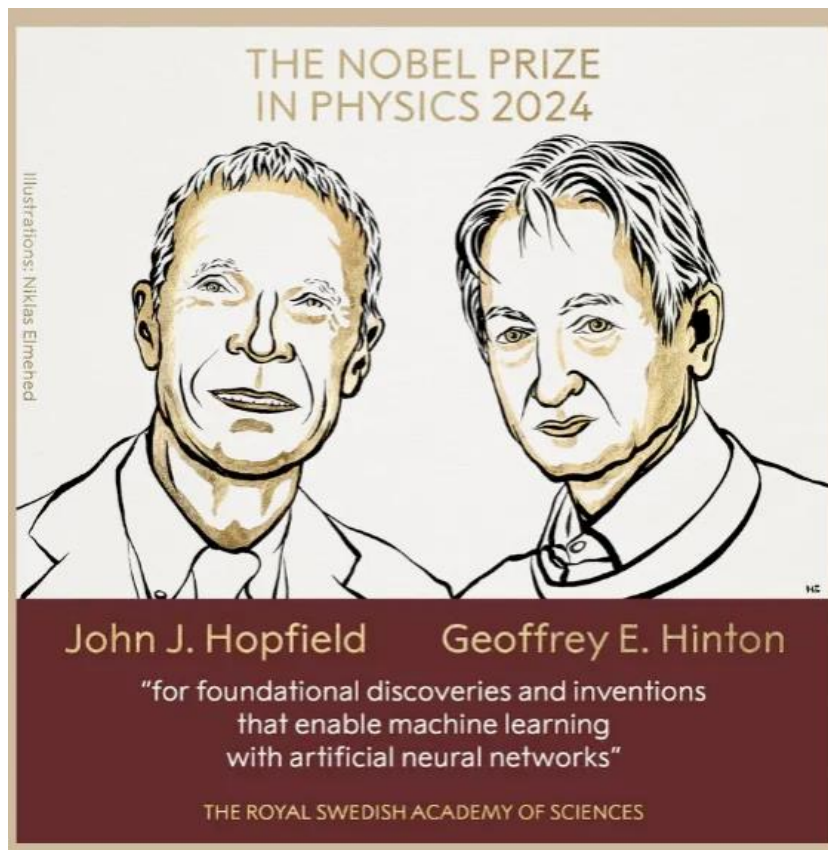
Shrnutí aktuálního dění pro nejhledanější dotazy

The screenshot shows the Seznam.cz homepage with a search bar at the top. The main content area features a news article titled "Ukrajinská vláda se otfásá. Odejde nejméně polovina členů, včetně Kuleby". A callout box highlights the "Krátké shrnutí" (Short summary) section, which is generated by AI. The summary text reads: "Ukrajinský ministr zahraničí Dmytro Kuleba rezignoval, což je součástí velké rošády v ukrajinské vládě, kterou naplánoval prezident Zelenský. Tato rošáda, která zahrnuje odchod více než poloviny členů kabinetu, včetně ministra strategického průmyslu, několika vicepremiérů a zástupce vedoucího prezidentské kanceláře, má proběhnout na začátku podzimu, který prezident Zelenský označil za extrémně důležitý pro Ukrajinu. Očekává se, že tato změna souvisí s novým obdobím výzev pro Ukrajinu, včetně silného ruského tlaku na frontě a snahy Ukrajiny získat západní zbraně pro úder hlouběji v ruském území." Below the summary is a "Celý článek" (Full article) button and engagement statistics: "Líbí se 123" and "Komentáře 68". To the right, there is a "Podobné články" (Similar articles) section with three items: "České noviny Lipavský i Fiala po útoku v Poltavě obhajují dodávky pomoci Ukrajině", "Médium Tři scénáře k porážce Ruska", and "Hlídací pes Slovensko chce svůj díl z čínského koláče: „Copak Evropa neporoučuje lidská práva?“". At the bottom right, there is an advertisement for "Alpen fest style Limburger" with a 30% discount, priced at 37.90 Kč. The bottom navigation bar includes icons for "Mapy", "Bazar", "Reality", "Auta", "Slovník", and "TV program".

Sumarizace článků

Nadpisy zahraničních článků v češtině





<https://o-seznam.cz/kariera/458967-ai-vyzkumnik-generativnich-jazykovych-modelu/>

