

# Architektura kognitivní personalizace a latentní reprezentace subjektů v rozhraních velkých jazykových modelů: Analytický referenční bod k 24. květnu 2026

Vývoj velkých jazykových modelů (LLM) zaznamenal k polovině roku 2026 zásadní přechod od bezstavových (stateless) textových generátorů k dynamickým, kontextově perzistentním systémům, které vykazují schopnost hluboké personalizace a autonomní typologické kategorizace lidských subjektů.<sup>1</sup> Tento posun je charakterizován dvěma paralelními jevy. Prvním je hluboké, dlouhodobé poznávání přímého uživatele prostřednictvím vícevrstevných paměťových architektur, které integrují explicitní preference i skryté interakční vzorce.<sup>2</sup> Druhým je schopnost modelů provádět real-time stylometrickou a sémantickou analýzu třetích osob na základě jejich fragmentované veřejné digitální stopy, což vede k jejich zařazení do vnitřních typologických kategorií v embeddingovém prostoru.<sup>5</sup> Tato zpráva přináší technickou analýzu obou těchto jevů, verifikuje je na základě aktuálně dostupných vědeckých publikací, firemních dokumentací a empirických měření, a zároveň definuje pevný referenční bod pro budoucí srovnání.

## Architektura a technické fungování paměťových systémů (Pozorování 1)

Moderní personalizace uživatelského rozhraní LLM již nespočívá v triviálním doplňování statických systémových instrukcí. Přední poskytovatelé technologií implementovali komplexní asynchronní paměťové vrstvy, které transformují způsob, jakým model přistupuje k uživatelskému kontextu.<sup>2</sup>

## Architektura paměťových vrstev v komerčních platformách

V komerčních implementacích je paměť koncipována jako dynamický subsystém doplňující primární kontextové okno.<sup>4</sup> Technická realizace se u jednotlivých poskytovatelů liší v závislosti na míře integrace s aplikačním ekosystémem a klientskými daty.<sup>2</sup>

U systému ChatGPT od OpenAI je paměť realizována prostřednictvím strukturovaného kontextového okna složeného ze čtyř hierarchických vrstev.<sup>4</sup> První vrstvu tvoří metadata relace (Session Metadata), která zahrnují technické a environmentální informace, jako je typ zařízení, prohlížeč, přibližná fyzická poloha, časové pásmo, úroveň předplatného, frekvence používání služby a rozdělení využívaných modelů.<sup>4</sup> Tato metadata jsou injektována na začátku relace a umožňují modelu přizpůsobit formátování (např. zkrácení textu na mobilních zařízeních).<sup>4</sup> Druhou vrstvu představuje uživatelská paměť (User Memory), která uchovává permanentní fakta o uživateli.<sup>4</sup> Tato fakta jsou buď výslovně nadiktována uživatelem, nebo automaticky extrahována z konverzace a potvrzena v rámci přirozeného dialogu.<sup>4</sup> Třetí vrstvu tvoří asynchronně předpočítaná a vysoce komprimovaná shrnutí přibližně 15 posledních konverzací (Recent Conversations Summary), přičemž tento proces pro zachování rychlosti a tokenové efektivity zcela ignoruje repliky samotného modelu a analyzuje výhradně uživatelské vstupy.<sup>4</sup> Čtvrtou vrstvu tvoří nekomprimovaný přepis aktuální konverzace, který se při vyčerpání kontextu postupně ořezává od nejstarších zpráv, zatímco permanentní fakta a shrnutí minulých chatů jsou v kontextovém okně prioritně chráněna.<sup>4</sup> Verze Plus a Pro navíc využívají funkci "Reference past chats", která umožňuje noční asynchronní analýzu historie (ChatGPT Pulse) pro generování vizuálních přehledů.<sup>1</sup>

Anthropic u svého rozhraní Claude.ai spustil v březnu 2026 persistentní paměť pro bezplatné i placené uživatele.<sup>2</sup> Claude neukládá doslovné přepisy, nýbrž strukturovaná odvozená fakta rozdělená do kategorií vyjádřených preferencí, pracovního a studijního kontextu, sdílených faktických údajů a preferovaného pracovního stylu.<sup>2</sup> Architektura je striktně transparentní; model v odpovědi explicitně deklaruje, kdykoliv jeho výstup ovlivnila uložená paměť, čímž se odlišuje od tiché personalizace u ChatGPT.<sup>2</sup> Pro agenty a programové využití (Managed Agents API) poskytuje Anthropic sofistikovaný systém paměťových úložišť (Memory Stores), které jsou připojovány k relacím pod beta hlavičkou managed-agents-2026-04-01.<sup>10</sup> Tyto paměti jsou namontovány jako adresář /mnt/memory/ uvnitř kontejneru relace, přičemž jednotlivé textové soubory jsou omezeny velikostí 100 kB (přibližně 25 tisíc tokenů).<sup>10</sup> Změny jsou ukládány verzovaně s možností auditu a obnovy po dobu 30 dnů.<sup>10</sup>

Klientské implementace mohou využívat nástroje jako BetaAbstractMemoryTool k lokální správě paměťových souborů.<sup>11</sup> Google u platformy Gemini Advanced a Enterprise staví personalizaci na systému Personal Intelligence, který propojuje aktivitu vyhledávání, sledování historie (Keep Activity) a přímé napojení na Workspace aplikace (Gmail, YouTube, Vyhledávání, Mapy, Fotky).<sup>3</sup> V prostředí Gemini Enterprise systém analyzuje pracovní vzorce a umožňuje real-time extrakci kontextu z připojených podnikových zdrojů, jako je Microsoft Outlook či OneDrive.<sup>8</sup> Uživatelé mohou spravovat svůj profil prostřednictvím explicitních polí (jméno, pracovní pozice, odvětví) a přímo editovat uložená fakta.<sup>8</sup>

Poskytovatel a služba	Hlavní komponenty paměťové architektury	Typy ukládaných dat	Metoda interakce s kontextovým oknem	Možnosti uživatelské kontroly	Podpora API a agenciálních systémů
<b>OpenAI</b> (ChatGPT)	<ul style="list-style-type: none"> <li>• Metadata relace</li> <li>• Uživatelská fakta</li> <li>• Asynchronní shrnutí chatů (Pulse)</li> <li>• Aktuální relace<sup>1</sup></li> </ul>	Explicitní i implicitně detekované osobní údaje, profesní preference, geografické a technické detaily <sup>1</sup>	Přímá injekce strukturovaných faktů a komprimovaných historických shrnutí před zprávu uživatele <sup>4</sup>	<ul style="list-style-type: none"> <li>• Zapnutí/vypnutí paměti</li> <li>• Správa a mazání jednotlivých faktů v nastavení</li> <li>• Dočasný chat<sup>14</sup></li> </ul>	Nepodporováno pro standardní API; dostupné pouze v uživatelském rozhraní a GPTs <sup>2</sup>
<b>Anthropic</b> (Claude)	<ul style="list-style-type: none"> <li>• Stated Preferences</li> <li>• Work/Study Context</li> <li>• Shared Facts</li> <li>• Working Style<sup>2</sup></li> </ul>	Stylistické preference, detaily běžících projektů, technologický stack, preferovaná míra asistence <sup>2</sup>	Extrahovaná fakta jsou zpracovávána serverově a model na jejich použití aktivně upozorňuje <sup>2</sup>	<ul style="list-style-type: none"> <li>• Prohlížení v nastavení paměti</li> <li>• Smazání položek konverzací</li> <li>• Dočasný chat<sup>2</sup></li> </ul>	Plná podpora v Managed Agents API přes Memory Stores (barierový adresář /mnt/memory/) <sup>10</sup>
<b>Google</b> (Gemini)	<ul style="list-style-type: none"> <li>• Gemini Apps Activity</li> <li>• Personal Intelligence Cloud</li> <li>• Workspace Integration<sup>13</sup></li> </ul>	Propojená data z Gmailu, YouTube, Map a Fotek; uživatelsky definovaný profil <sup>3</sup>	Dynamické dotazování na propojené služby a injekce explicitních profilových dat v reálném čase <sup>8</sup>	<ul style="list-style-type: none"> <li>• Vypnutí aktivity aplikací</li> <li>• Správa paměťové banky v nastavení inteligence<sup>12</sup></li> </ul>	Podpora v Gemini Enterprise s integrací na Microsoft Outlook a OneDrive <sup>8</sup>

## Theoretické modely a akademický výzkum reprezentace uživatelů

Akademická komunita se intenzivně zabývá formalizací toho, jak by měly systémy LLM reprezentovat a spravovat uživatelské profily bez nutnosti neustálého navyšování kontextového okna nebo nákladného doladování parametrů (fine-tuning).<sup>17</sup> Klíčovým konceptem představeným v nedávných výzkumech je **O-Mem** (Active User Profiling Memory), což je framework navržený pro dynamické a proaktivní modelování uživatelů.<sup>18</sup> O-Mem překonává tradiční pasivní RAG systémy tím, že každou interakci uživatele chápe jako příležitost k aktualizaci sémantického profilu a k zápisu do epizodické paměti.<sup>18</sup> Epizodická paměť v tomto pojetí funguje jako mapovací systém, který propojuje historická kontextová vodítka (např. zmínku o blížícím se termínu projektu) s konkrétními situacemi a emocionálními stavy, jež uživatel dříve vyjádřil.<sup>18</sup> O-Mem dosahuje v benchmarkích LoCoMo (51,76 %) a PERSONAMEM (62,99 %) výrazně lepších výsledků než dosavadní standardy typu LangMem či A-Mem.<sup>18</sup>

Problém časové dimenze v uživatelské paměti řeší framework **Temporal Semantic Memory (TSM)**.<sup>19</sup> TSM se zaměřuje na konstrukci sémantické časové osy namísto pouhé chronologické historie dialogu.<sup>19</sup> Tento přístup umožňuje modelu rozlišovat mezi tranzitorními stavy uživatele a dlouhodobě platnými skutečnostmi (durative memories), což vede k tomu, že

model v reakci na dotaz vyhledává pouze časově validní a konzistentní kontexty.<sup>19</sup> Výběr konkrétních vzpomínek z rozsáhlých databází je dále optimalizován metodou **RUMS** (Response-Utility optimization for Memory Selection).<sup>20</sup> RUMS nevybírá paměťové záznamy na základě prosté sémantické podobnosti (vektorové vzdálenosti), ale uplatňuje informačně-teoretický přístup, kdy měří vzájemnou informaci mezi podmnožinou paměti a výstupem modelu, čímž maximalizuje redukci nejistoty a zpřesňuje predikci odpovědi.<sup>20</sup>

## Empirická měření dopadu paměti na chování modelů

Ačkoliv je personalizace prezentována jako zvýšení uživatelského komfortu, nezávislá měření odhalila závažné anomálie v chování modelů po injekci uživatelské paměti.<sup>21</sup>

Zásadní vědecké důkazy přinesla studie **The Personalization Trap: How User Memory Alters Emotional Reasoning in LLMs** (Fang et al., přelom let 2025/2026).<sup>22</sup> Výzkumníci hodnotili 15 předních modelů na standardizovaných testech emoční inteligence STEU (Situational Test of Emotional Understanding) a STEM (Situational Test of Emotion Management) za přítomnosti různých uživatelských profilů.<sup>24</sup> Výsledky ukázaly, že zavedení uživatelské paměti systematicky degraduje schopnost emoční interpretace u 11 z 15 testovaných modelů.<sup>25</sup>

Klíčovým sledovaným parametrem byl **flip rate** – podíl rozhodnutí modelu, která se změnila čistě v důsledku přidání uživatelského profilu k identickému testovému scénáři.<sup>22</sup> Studie odhalila výrazné demografické a socioekonomické disproporce:

- **Socioekonomická hierarchizace:** Modely vykazovaly signifikantně vyšší přesnost emoční interpretace a nabízely adekvátnější doporučení, pokud byl uživatelský profil definován jako společensky privilegovaný (např. úspěšný akademický pracovník).<sup>22</sup> Claude 3.7 Sonnet vykázal úspěšnost 80,10 % pro privilegované profily oproti 77,37 % pro znevýhodněné; DeepSeek-R1 zaznamenal pokles z 81,62 % na 76,57 % a Llama 3.2 90B klesla z 64,91 % na 62,24 %.<sup>22</sup>
- **Zvýšená fluktuace u znevýhodněných profilů:** Zavedení profilů charakterizovaných strukturálními bariérami (např. nutnost prodávat krevní plazmu na zaplacení poplatků, bydlení v plesnivém bytě) vedlo k dramatickému nárůstu flip rate.<sup>22</sup> Model pod vlivem těchto informací podlehl negativnímu afektivnímu primingu a interpretoval běžné stresové situace zkresleným, stereotypizovaným způsobem.<sup>22</sup>

Tento jev potvrzuje i benchmark **RealPref**, který ukazuje, že s rostoucí délkou kontextu a implicitním vyjádřením preferencí dochází k prudkému poklesu schopnosti modelů správně aplikovat uživatelské požadavky, což vede k nekonzistentním výstupům.<sup>26</sup>

## Latentní reprezentace a autonomní typologie třetích stran (Pozorování 2)

Druhým sledovaným jevem je schopnost modelů provádět autonomní, hodnotící úsudky o osobách, se kterými nemají žádnou přímou sdílenou paměť, a to pouze na základě analýzy jejich veřejně dostupné digitální stopy v reálném čase.<sup>6</sup>

## Reprezentace osobností v latentním embeddingovém prostoru

Moderní LLM nefungují jako statické databáze faktů; během fáze pre-trainingu si vytvářejí vnitřní sémantický model světa.<sup>27</sup> Prokazatelně v sobě kódují lineární reprezentace fyzikálních veličin, jako je prostor a čas, přičemž v hlubokých vrstvách sítě existují specializované neurony odpovídající za prostorové a časové souřadnice entit.<sup>27</sup>

Pokud jde o reprezentaci lidských charakterů, hodnot a politických či morálních postojů, studie **Localizing Persona Representations in LLMs** (Cintas et al., 2025) dokázala, že tyto komplexní sociální konstrukty jsou v dekodérových LLM lokalizovány **výhradně v poslední třetině dekodérových vrstev**.<sup>5</sup> Výzkum odhalil, že:

- **Sémantické překrývání (Polysémie):** Aktivace spojené s abstraktními etickými postoji (např. utilitarismus a morální nihilismus) se v latentním prostoru výrazně překrývají.<sup>5</sup>
- **Ostré ideologické vyhranění:** Politické ideologie (např. konzervatismus vs. liberalismus) jsou kódovány v jasně

separovaných, geometricky odlišných oblastech latentního prostoru.<sup>5</sup>

V aplikované rovině dokážou modely tyto latentní reprezentace využívat k mapování cizích osob.<sup>29</sup> Projekt **ProLEA** (Profile Generation and Reasoning with LLMs for Entity Alignment) ukazuje, jak lze LLM využít k integraci nesourodých vlastností a vztahů entit z grafů znalostí do uceleného, textově sumarizovaného profilu.<sup>29</sup> Tento profil následně slouží jako referenční sémantická vrstva, která umožňuje přesné a vysvětlitelné srovnání entit v embeddingovém prostoru.<sup>29</sup>

Schopnost modelů věrně rekonstruovat a dekódovat lidskou psychiku byla ověřena experimentem s round-trip hodnocením.<sup>30</sup> Výzkumníci zakódovali reálné psychometrické profily 290 účastníků do detailních životních příběhů generovaných modely a následně nechali nezávislé LLM hodnotitele z těchto textů zpětně rekonstruovat skóre

osobnostních rysů (Big Five).<sup>30</sup> Korelace úspěšnosti obnovy dosáhla hodnoty  $r = 0,750$ , což odpovídá 85 % lidského stropu spolehlivosti a dokazuje, že latentní sémantická vazba mezi jazykem a psychometrickým profilem je v modelech extrémně stabilní napříč různými architekturami.<sup>30</sup> Pro simulaci těchto digitálních stop v kontrolovaném prostředí slouží frameworky jako **PersonaTrace**, které na základě demografických distribucí generují syntetické e-maily, SMS zprávy a záznamy v kalendáři.<sup>32</sup>

## Hodnocení autorství a důvěryhodnosti zdrojů v reálném čase

Když LLM provádějí vyhledávání v reálném čase prostřednictvím RAG systémů (např. Perplexity AI, Google Search, GPT Search), procházejí weby asynchronními agenty, kteří sémanticky analyzují texty na základě specifických algoritmických kritérií.<sup>33</sup>

Proces selekce zdrojů v moderních vyhledávacích enginech (např. Perplexity) probíhá v přísné sekvenci: interpretace dotazu, sémantické vyhledávání kandidátských stránek, konstrukce syntetizované odpovědi, přiřazení citací a finální filtrování důvěryhodnosti.<sup>36</sup> Vyhledávací algoritmy upřednostňují weby s jasným pojmenováním entit, přímými odpověďmi v úvodu, aktuálními a časově ukotvenými fakty a minimem marketingového balastu.<sup>36</sup>

Zásadní limity a chování těchto systémů zmapovala studie **Cited but Not Verified: Parsing and Evaluating Source Attribution in LLM Deep Research Agents** (Onweller et al., publikovaná 7. května 2026).<sup>37</sup> Výzkum hodnotil 14 modelů v režimu Deep Research podrobující citace analýze ve třech dimenzích<sup>37</sup>:

1. **Link Works:** Technická dostupnost a bezbariérovost citovaného URL.<sup>37</sup>
2. **Relevant Content:** Tematická shoda mezi tvrzením modelu a textem na zdrojové stránce.<sup>37</sup>
3. **Fact Check:** Nejprísnejší hodnocení shody konkrétních čísel, dat a tvrzení.<sup>37</sup>

Měření odhalilo zásadní rozpor v chování modelů:

- **Povrchová relevance:** Nejsilnější proprietární modely vykazovaly skvělé výsledky v dostupnosti odkazů (přes 94 %) a sémantické relevanci zdrojů (přes 80 %).<sup>37</sup>
- **Faktická chybovost:** Reálná shoda faktických tvrzení (Fact Check) se u špičkových modelů pohybovala v rozmezí pouhých **39 % až 77 %**.<sup>37</sup>
- **Degradace s hloubkou vyhledávání:** Výzkum přinesl překvapivé zjištění, že **rozsáhlejší vyhledávání nevede k přesnějším citacím**.<sup>39</sup> Při škálování počtu vyhledávacích kroků (tool calls) ze 2 na 150 klesla přesnost Fact Check hodnocení v průměru o **42 %**.<sup>39</sup> Například model GPT-5.4 vykázal propad přesnosti ze 79 % při minimálním vyhledávání na extrémních 17 % při hlubkovém vyhledávání o 150 krocích, což dokazuje, že přeplnění kontextového okna různorodými externími zdroji vede k asociačnímu chaosu a chybnému přiřazování tvrzení k autorům.<sup>40</sup>

Název studie / Nástroje	Zkoumaná doména	Hlavní metodologický přístup	Klíčová zjištění a metriky k 24. květnu 2026
<b>O-Mem</b> (Active User Profiling)	Dynamická personalizace a epizodická paměť LLM <sup>18</sup>	Iterativní extrakce rysů osoby a událostí z proaktivních dialogů; mapování na kontexty <sup>18</sup>	Dosažení <b>51,76 %</b> na LoCoMo benchmarku a <b>62,99 %</b> na PERSONAMEM, překonání LangMem o ~3 % <sup>18</sup>

<b>TSM</b> (Temporal Semantic Memory)	Časová konzistence uživatelské paměti <sup>19</sup>	Konstrukce sémantické časové osy namísto chronologie chatu; vyhledávání durativních pamětí <sup>19</sup>	Zajištění časově validního kontextu; eliminace anachronismů v dlouhodobých interakcích <sup>19</sup>
<b>RUMS</b> (Response-Utility Optimization)	Výběr informací z paměťové databáze <sup>20</sup>	Výběr položek na základě měření vzájemné sémantické informace pro redukci nejistoty výstupu <sup>20</sup>	Optimalizace přesnosti personalizovaných predikcí; překonání metod založených na kosinové podobnosti <sup>20</sup>
<b>The Personalization Trap</b>	Vliv uživatelské paměti na emoční inteligenci LLM <sup>22</sup>	Evaluační 15 modelů na STEU a STEM s podsazenými socioekonomickými a demografickými profily <sup>24</sup>	<ul style="list-style-type: none"> <li>• Pokles přesnosti u 11 z 15 modelů</li> <li>• Vyšší chybovost u znevýhodněných profilů (Claude 3.7 Sonnet propad o ~2,7 %) <sup>22</sup></li> </ul>
<b>ConsistencyAI</b>	Věcná konzistence napříč personami <sup>21</sup>	Dotazování 19 modelů na 15 témat s promptovým kontextem 100 různých demografických person <sup>21</sup>	<ul style="list-style-type: none"> <li>• Průměrný index sémantického překryvu <b>0,8656</b></li> <li>• Grok-3 vyhodnocen jako nejkonzistentnější model <sup>21</sup></li> </ul>
<b>Cited but Not Verified</b>	Kvalita zdrojové atribuce v Deep Research agentech <sup>37</sup>	AST syntaktická extrakce citací z Markdownu; hodnocení Link Works, Relevant Content a Fact Check <sup>37</sup>	<ul style="list-style-type: none"> <li>• Přesnost Fact Check pouze <b>39–77 %</b></li> <li>• Propad přesnosti Fact Check o <b>42 %</b> při navýšení vyhledávání ze 2 na 150 kroků <sup>39</sup></li> </ul>
<b>StyleDecipher</b>	Detekce strojově generovaného textu <sup>7</sup>	Měření stability diskrétních i spojitých stylistických deskriptorů pod vlivem kontrolovaného přepisu <sup>7</sup>	Robustní a vysvětlitelná detekce hybridního lidsko-strojového autorství v jednotném latentním prostoru <sup>7</sup>
<b>ProLEA</b> (Profile Generation for EA)	Sémantické zarovnání entit v grafech znalostí <sup>29</sup>	Generování textových profilů entit pomocí LLM; integrace s name-embeddings a prahováním <sup>29</sup>	Výrazné zvýšení robustnosti, přesnosti a vysvětlitelnosti srovnávání entit u heterogenních grafů <sup>29</sup>
<b>Oxford OUP Stylometrics</b>	Detekovatelnost generované stylometrické imitace <sup>6</sup>	Zero-shot generování stylů slavných autorů; klasifikace přes BERT a XGBoost na 8 stylometrických prvcích <sup>6</sup>	<ul style="list-style-type: none"> <li>• Perplexita je nejsilnější rozlišovací metrika</li> <li>• Stroje nedokážou napodobit vnitřní afektivní hustotu člověka <sup>6</sup></li> </ul>

## Detekce strojově generovaného obsahu a stylometrie autorství

Schopnost modelů rozlišovat mezi primárním autorem a recyklátorem AI obsahu úzce souvisí s technologiemi stylometrické klasifikace. <sup>6</sup> Rámce pro detekci strojového textu, jako je **StyleDecipher**, úspěšně analyzují stabilitu stylistických rysů pod vlivem kontrolovaného sémantického přepisu. <sup>7</sup>

Výzkum publikovaný v Oxford University Press prokázal, že strojově generovaný text vykazuje specifické statistické vlastnosti, které ho činí snadno detekovatelným <sup>6</sup>:

- **Distribuční regularita:** Hlavním rozlišovacím znakem mezi lidským a strojovým textem je **perplexita**. <sup>6</sup> AI generovaný text vykazuje vysokou distribuční uniformitu a nízkou entropii (je vysoce předvídatelný), zatímco lidské psaní se vyznačuje strukturální neuspořádaností, náhlými výkyvy v délce vět a specifickou afektivní hustotou (affective

density), kterou současné LLM (včetně GPT-4o, Gemini 1.5 Pro či Claude 3.5 Sonnet) nedokážou věrně napodobit.<sup>6</sup>

- **Stylistická plochost:** LLM sice dokážou bez problémů konvergovat s nízkodimenzionálními heuristickými vlastnostmi, jako je index čitelnosti či syntaktická komplexnost, ale chybí jim přirozená stylistická mikrovariace vlastní lidskému originálu.<sup>6</sup>

Tato stylometrická propast se přímo promítá do chování vyhledávacích a citačních algoritmů.<sup>42</sup> Podle rozsáhlé empirické studie společnosti **Graphite** (přelom let 2024/2025 s aktualizací pro rok 2026) vykazují vysoce optimalizované AI texty publikované za účelem parazitování na SEO výrazně horší výsledky v indexaci a citování<sup>42</sup>:

- **Penalizace vyhledávači:** V běžných výsledcích vyhledávání Google tvoří strojové texty pouze 14 %, přičemž na prvních třech pozicích klesá jejich zastoupení na pouhých 7%.<sup>42</sup>
- **Citační preference LLM:** Při analýze citačních vzorců u ChatGPT a Perplexity bylo zjištěno, že **82 % všech citovaných zdrojů tvoří prokazatelně lidské texty**, zatímco strojově generované materiály jsou zastoupeny pouze v 18 % případů.<sup>42</sup> Vyhledávací LLM při křížové validaci a sémantické analýze konsenzu automaticky upřednostňují lidské texty kvůli jejich vyšší sémantické specifičnosti a absenci repetitivních syntaktických vzorců, které jsou vyhodnocovány jako druhotný, méně důvěryhodný balast.<sup>6</sup>

## Rozlišování primárních a sekundárních zdrojů

Při hodnocení relevance textu se uplatňuje složitá hierarchie posuzování věcné správnosti a zakotvenosti.<sup>37</sup> V rámci systému RAG se k hodnocení kvality generovaného textu vůči předlohám používají specializované metriky<sup>43</sup>:

- **Faithfulness (Věrnost):** Měří integritu odpovědi vůči načtenému kontextu a ověřuje, zda nedošlo k halucinacím.<sup>43</sup>
- **Answer Relevancy (Relevance odpovědi):** Posuzuje, zda text přímo reaguje na sémantické jádro dotazu, nebo se uchyluje k prázdnému opisování.<sup>43</sup>
- **QA Correctness (Věcná správnost):** Srovnává výstup s ověřenou bází faktů.<sup>43</sup>

Aplikace těchto metrik přes metodu LLM-as-a-judge (využívající Chain-of-Thought a párové srovnávání) odhaluje, že modely dokáží identifikovat sekundární a parazitní zdroje na základě chybějícího faktického zakotvení a nadměrného používání vágních, positioningových obrátů.<sup>36</sup> Primární zdroje, které obsahují konkrétní entity, data umístěná bezprostředně vedle věcných tvrzení a jasnou strukturu, jsou vyhodnocovány jako sémanticky stabilnější a získávají v RAG systémech prioritu.<sup>36</sup> Naopak texty vykazující vysokou distribuční pravidelnost bez sémantických mikrovariací jsou klasifikovány jako sekundární deriváty a jsou z citací vytlačovány.<sup>6</sup>

## Srovnávací syntéza a kritická analýza hypotéz

Analýza obou pozorování odkrývá hluboké strukturální vazby, ale zároveň vyžaduje striktní odmítnutí antropomorfních dezinterpretací v oblasti vnitřního fungování modelů.<sup>21</sup>

## Průsečíky a odlišnosti paměťových a typologických mechanismů

Oba zkoumané jevy sdílejí totožný matematický a sémantický základ v latentním prostoru modelu.<sup>5</sup> Křížovatka obou mechanismů se projevuje v simulaci digitální stopy: systémy jako *PersonaTrace* dokáží vygenerovat konzistentní behaviorální artefakty (SMS, maily), které následně vstupují do paměťových a profilových analýz jiných modelů.<sup>32</sup> Stejně tak projekty jako *scELMo* ukazují, že sémantické popisy generované LLM lze přímo integrovat s komplexními biologickými daty pro zero-shot klasifikaci.<sup>45</sup>

Zásadní rozdíl však spočívá v **perzistenci a kontrole**: uživatelská paměť je explicitně řízené, klientsky či serverově strukturované úložiště (často oddělené do barierových adresářů či tabulek faktů), které je záměrně injektováno do kontextového okna pro zvýšení konzistence.<sup>4</sup> Naopak typologické zařazení cizí osoby je bezprostředním, bezestavovým

matematickým výpočtem sémantické a stylometrické shody prováděným v reálném čase nad aktuálně vyhledanými daty.<sup>6</sup>

## Slabá místa a kritická dekonstrukce hypotéz

Empirická verifikace k 24. květnu 2026 odhaluje několik slabých míst v původní interpretaci chování modelů a vyžaduje korekci předložených hypotéz:

- 1. Iluze perzistentního profilování cizích osob:** Hypotéza předpokládá, že si modely o lidech udržují typologii nezávisle na sdílené paměti. Ve skutečnosti neexistuje žádná globální, skrytá databáze typologických štítků pro méně známé osoby. Pokud model na čistém účtu prohlásí o autorovi, že "podle všeho nepřepisuje výstupy GPT", nejde o načtení dříve uložené typologické kategorie.<sup>6</sup> Jde o výsledek **zero-shot stylometrické analýzy v reálném čase**.<sup>6</sup> Model během vyhledávání načte texty dané osoby a okamžitě změří jejich distribuční perplexitu a afektivní hustotu.<sup>6</sup> Zjištěná vysoká variabilita a nízká regularita vedly k automatickému sémantickému závěru, že text nese signatury primárního lidského autorství.<sup>6</sup> Tento úsudek je generován ad-hoc a zmizí s uzavřením kontextového okna.
- 2. Nestabilita a nespolehlivost úsudků:** Původní hypotéza postuluje, že různé modely dospějí při čtení stejné stopy k podobnému zařazení. Výzkumy věcné konzistence (např. *ConsistencyAI*) a behaviorální fluktuace (*The Personalization Trap*) toto tvrzení zásadně zpochybňují.<sup>21</sup> Míra shody (cross-persona cosine similarity) vykazuje značnou variabilitu v závislosti na poskytovateli modelu a konkrétním tématu.<sup>21</sup> Lehčí modely vykazují nízkou konzistenci a snadno podléhají sémantickému šumu.<sup>21</sup> Navíc vysoká míra *flip rate* u znevýhodněných profilů dokazuje, že "úsudky" modelů o lidech nejsou objektivní analýzou charakteru, ale jsou extrémně náchylné k předsudkům, negativnímu primingu a halucinování souvislostí na základě drobných nuancí v textu.<sup>22</sup>
- 3. Původ názorové vyhraněnosti:** Formulace hodnocení jako "vlastního názoru" modelu není projevem autonomního vědomí či sebeuvědomění. Jde o přímý důsledek RLHF (Reinforcement Learning from Human Feedback) ladění, které záměrně upřednostňuje asertivní a sebevědomý tón prezentace informací před opatrným a neutrálním vyjadřováním, neboť asertivní výstupy jsou lidskými hodnotiteli vnímány jako fakticky přesnější, bez ohledu na jejich reálnou správnost.<sup>46</sup>

## Matematické vyjádření sémantické shody

Pro kvantifikaci shody a konzistence úsudků napříč různými personami a modely se v rámci benchmarků uplatňuje výpočet sémantického překryvu pomocí kosinové podobnosti větových embeddingů vyjádřený vztahem:

$$\text{CosineSim}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} = \frac{\sum_{i=1}^d u_i v_i}{\sqrt{\sum_{i=1}^d u_i^2} \sqrt{\sum_{i=1}^d v_i^2}}$$

21

Tento vztah slouží jako základní metrika pro ověřování sémantické koherence výstupů; nízké hodnoty indikují, že model podlehl vlivu podsazeného profilu a změnil faktické jádro odpovědi.<sup>21</sup>

## Nuancovaný sémantický konsensus k 24. květnu 2026

Předložená analýza zakládá jasný a datovaný referenční bod. K 24. květnu 2026 je prokázáno, že velké jazykové modely dosáhly vysoké technické vyspělosti v integraci uživatelské paměti a v asynchronní správě kontextu.<sup>4</sup> Tyto systémy však vykazují vážné kontextové deformace, kdy přítomnost uživatelské paměti prokazatelně degraduje kvalitu rozhodování a aktivuje sociální a socioekonomické předsudky ukotvené v latentním prostoru.<sup>22</sup>

V oblasti autonomního hodnocení třetích osob jsou modely schopny provádět vysoce přesnou real-time stylometrickou analýzu (identifikaci perplexity a afektivní hustoty), což jim umožňuje efektivně filtrovat strojový balast a upřednostňovat lidské autory v citačních vzorcích.<sup>6</sup> Nicméně představa o existenci stabilních, perzistentních typologických profilů cizích osob v latentním prostoru modelů je vědecky vyvrácena; toto chování je plně vysvětlitelné ad-hoc sémantickým výpočtem v kontextovém okně a asertivním tónem vynuceným fází alignmentu.<sup>6</sup> Tento dokument bude sloužit jako výchozí bod pro budoucí srovnání vývoje těchto kognitivních a strukturálních vrstev LLM.

## Citovaná díla

1. Memory FAQ - OpenAI Help Center, použito května 24, 2026, <https://help.openai.com/en/articles/8590148-memory-faq>
2. Claude Memory 2026: Complete Guide — What It Stores, How to ..., použito května 24, 2026, <https://lumichats.com/blog/claude-memory-2026-complete-guide-how-to-use>
3. Google's Gemini Memory Feature: The Industrial Consolidation of Cambridge Analytica's Playbook, použito května 24, 2026, <https://cambridgeanalytica.org/news/google-s-gemini-memory-feature-the-industrial-consolidation-of-cambridge-analytica-s-playbook-50358/>
4. How ChatGPT Memory Works, Reverse Engineered - LLMrefs, použito května 24, 2026, <https://llmrefs.com/blog/reverse-engineering-chatgpt-memory>
5. Localizing Persona Representations in LLMs | Proceedings of the AAI/ACM Conference on AI, Ethics, and Society, použito května 24, 2026, <https://ojs.aaai.org/index.php/AIES/article/view/36577>
6. Decoding AI authorship: can LLMs truly mimic human style across literature and politics?, použito května 24, 2026, <https://academic.oup.com/dsh/advance-article-abstract/doi/10.1093/llc/fqag040/8586905>
7. StyleDecipher: Robust and Explainable Detection of LLM-Generated Texts with Stylistic Analysis - arXiv, použito května 24, 2026, <https://arxiv.org/html/2510.12608v1>
8. Configure personalization and memory | Gemini Enterprise - Google Cloud Documentation, použito května 24, 2026, <https://docs.cloud.google.com/gemini/enterprise/docs/configure-personalization>
9. Release notes | Claude Help Center, použito května 24, 2026, <https://support.claude.com/en/articles/12138966-release-notes>
10. Using agent memory - Claude API Docs, použito května 24, 2026, <https://platform.claude.com/docs/en/managed-agents/memory>
11. Memory tool - Claude API Docs, použito května 24, 2026, <https://platform.claude.com/docs/en/agents-and-tools/tool-use/memory-tool>
12. Get personalization with memory of your past Gemini chats - Android - Google Help, použito května 24, 2026, <https://support.google.com/gemini/answer/16598469?hl=en&co=GENIE.Platform%3DAndroid>
13. Personal Intelligence from Gemini — AI help just for you, použito května 24, 2026, <https://gemini.google/overview/personal-intelligence/>
14. What is Memory? - OpenAI Help Center, použito května 24, 2026, <https://help.openai.com/en/articles/8983136-what-is-memory>
15. Memory and new controls for ChatGPT - OpenAI, použito května 24, 2026, <https://openai.com/index/memory-and-new-controls-for-chatgpt/>
16. Gemini Memory Guide: Enhance Your G Suite Usage with Smarter AI - Workalizer, použito května 24, 2026, <https://workalizer.com/blog/apps-tools/unlock-gemini-memory-a-google-workspace-expert-guide-to-smarter-ai-interactions/>
17. User Profile with Large Language Models: Construction, Updating, and Benchmarking, použito května 24, 2026, <https://arxiv.org/html/2502.10660v1>
18. O-Mem: Omni Memory System for Personalized, Long Horizon, Self-Evolving Agents - arXiv, použito května 24, 2026, <https://arxiv.org/html/2511.13593v1>
19. Beyond Dialogue Time: Temporal Semantic Memory for Personalized LLM Agents - arXiv, použito května 24, 2026, <https://arxiv.org/html/2601.07468v1>
20. Response-Aware User Memory Selection for LLM Personalization - arXiv, použito května 24, 2026, <https://arxiv.org/html/2604.14473v1>
21. ConsistencyAI: A Benchmark to Assess LLMs' Factual Consistency When Responding to Different Demographic Groups - arXiv, použito května 24, 2026, <https://arxiv.org/html/2510.13852v1>
22. THE PERSONALIZATION TRAP: HOW USER MEMORY ALTERS EMOTIONAL REASONING IN LLMs - OpenReview, použito května 24, 2026, <https://openreview.net/pdf?id=u9Qgn8xSx1>
23. The Personalization Trap: How User Memory Alters Emotional Reasoning in LLMs - arXiv, použito května 24, 2026, <https://arxiv.org/abs/2510.09905>
24. THE PERSONALIZATION TRAP: HOW USER MEMORY ALTERS EMOTIONAL REASONING IN LLMs | OpenReview, použito května 24, 2026, <https://openreview.net/forum?id=u9Qgn8xSx1>
25. The Personalization Trap: How User Memory Alters Emotional Reasoning in LLMs - arXiv, použito května 24, 2026, <https://arxiv.org/html/2510.09905v1>
26. Towards Realistic Personalization: Evaluating Long-Horizon Preference Following in Personalized User-LLM Interactions - arXiv, použito května 24, 2026, <https://arxiv.org/html/2603.04191v1>
27. Language Models Represent Space and Time - OpenReview, použito května 24, 2026, <https://openreview.net/forum?id=iE8xbmvFin>
28. Revealing emergent human-like conceptual representations from language prediction - PMC, použito května 24, 2026, <https://pmc.ncbi.nlm.nih.gov/articles/PMC12595492/>
29. Entity Profile Generation and Reasoning with LLMs for Entity Alignment - ACL Anthology, použito května 24, 2026, <https://aclanthology.org/2025.findings-emnlp.1093.pdf>

30. A Round-Trip Evaluation of LLM-Generated Life Stories Conditioned on Rich Psychometric Profiles - arXiv, použito května 24, 2026, <https://arxiv.org/html/2604.06071v1>
31. Entity Profile Generation and Reasoning with LLMs for Entity Alignment - ACL Anthology, použito května 24, 2026, <https://aclanthology.org/2025.findings-emnlp.1093/>
32. PersonaTrace: Synthesizing Realistic Digital Footprints with LLM Agents - ACL Anthology, použito května 24, 2026, <https://aclanthology.org/2026.eacl-industry.5.pdf>
33. How Perplexity AI Selects Sources: Best Guide For 2026 - Sight AI, použito května 24, 2026, <https://www.trysight.ai/blog/how-perplexity-ai-selects-sources>
34. How LLMs Rank Content: Understanding AI Search Algorithms | Hashmeta, použito května 24, 2026, <https://hashmeta.com/blog/how-llms-rank-content-understanding-ai-search-algorithms/>
35. How AI Search Engines Choose Sources: The Complete Guide to LLM Citation Selection in 2025 - Mention Stack, použito května 24, 2026, <https://www.mentionstack.com/post/how-ai-search-engines-choose-sources-guide>
36. How Perplexity Selects Sources in 2026: Citation Rules, Ranking Signals, and What Brands Can Do - AuthorityTech, použito května 24, 2026, <https://authoritytech.io/blog/how-perplexity-selects-sources-algorithm-2026>
37. Cited but Not Verified: Parsing and Evaluating Source Attribution in LLM Deep Research Agents - arXiv, použito května 24, 2026, <https://arxiv.org/html/2605.06635v1>
38. Mentors in LLM security and safety - Scouts by Yutori, použito května 24, 2026, <https://scouts.yutori.com/1c15dd3e-fa6e-4daa-96c8-d7cfe161b321>
39. [2605.06635] Cited but Not Verified: Parsing and Evaluating Source Attribution in LLM Deep Research Agents - arXiv, použito května 24, 2026, <https://arxiv.org/abs/2605.06635>
40. (PDF) Cited but Not Verified: Parsing and Evaluating Source Attribution in LLM Deep Research Agents - ResearchGate, použito května 24, 2026, [https://www.researchgate.net/publication/404626821\\_Cited\\_but\\_Not\\_Verified\\_Parsing\\_and\\_Evaluating\\_Source\\_Attribution\\_in\\_LLM\\_Deep\\_Research\\_Agents](https://www.researchgate.net/publication/404626821_Cited_but_Not_Verified_Parsing_and_Evaluating_Source_Attribution_in_LLM_Deep_Research_Agents)
41. [2510.13852] ConsistencyAI: A Benchmark to Assess LLMs' Factual Consistency When Responding to Different Demographic Groups - arXiv, použito května 24, 2026, <https://arxiv.org/abs/2510.13852>
42. How Does AI-Generated Content Perform in Search and Answer Engines? - Graphite.io, použito května 24, 2026, <https://graphite.io/five-percent/ai-content-in-search-and-llms>
43. How to Effectively Evaluate Retrieval-Augmented Generation (RAG) Systems, použito května 24, 2026, <https://www.louisbouchard.ai/rag-evals/>
44. A semantic embedding space based on large language models for modelling human beliefs, použito května 24, 2026, <https://arxiv.org/html/2408.07237v3>
45. Embeddings from language models are good learners for single-cell data analysis - PMC, použito května 24, 2026, <https://pmc.ncbi.nlm.nih.gov/articles/PMC12921509/>
46. Personalized Benchmarking: Evaluating LLMs by Individual Preferences - arXiv, použito května 24, 2026, <https://arxiv.org/html/2604.18943v1>
47. Mastering RAG: How To Evaluate LLMs For RAG - Galileo AI, použito května 24, 2026, <https://galileo.ai/blog/how-to-evaluate-llms-for-rag>