

# In-Browser Agents for Search Assistance

Saber Zerhoudi  
University of Passau  
Passau, Germany  
saber.zerhoudi@uni-passau.de

Michael Granitzer  
University of Passau  
Passau, Germany  
Interdisciplinary Transformation University Austria  
Linz, Austria  
michael.granitzer@uni-passau.de

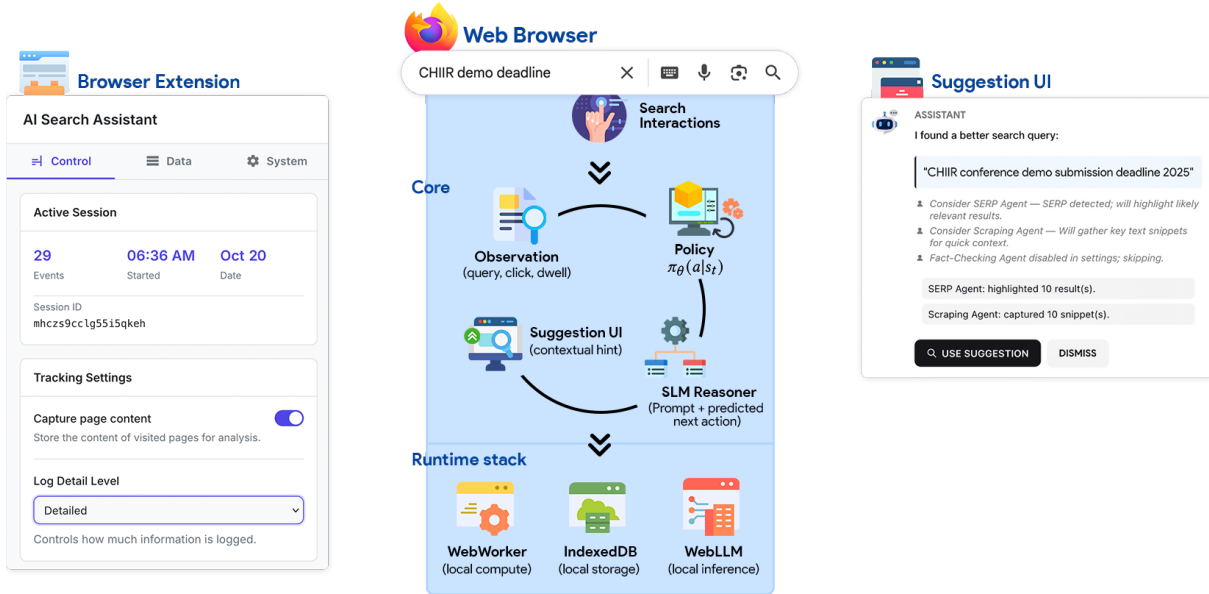


Figure 1: Overview and Interfaces of the In-Browser, Behavior-Grounded Search Assistant.

## Abstract

A fundamental tension exists between the demand for sophisticated AI assistance in web search and the need for user data privacy. Current centralized models require users to transmit sensitive browsing data to external services, which limits user control. In this paper, we present a browser extension<sup>1</sup> that provides a viable in-browser alternative. We introduce a hybrid architecture that functions entirely on the client side, combining two components: (1) an adaptive probabilistic model that learns a user’s behavioral policy from direct feedback, and (2) a Small Language Model (SLM), running in the browser, which is grounded by the probabilistic model to generate context-aware suggestions. To evaluate this approach, we conducted a three-week longitudinal user study with 18 participants. Our results show that this privacy-preserving approach is highly effective at adapting to individual user behavior, leading to measurably improved search efficiency. This work demonstrates that sophisticated AI assistance is achievable without compromising user privacy or data control.

## CCS Concepts

• **Information systems** → **Personalization**; **Web search engines**; *Users and interactive retrieval*; • **Security and privacy** → *Privacy-preserving technologies*; • **Human-centered computing** → *User models*.

## Keywords

Search Personalization, User Modeling, Browser Extension, Small Language Models

## 1 Introduction

The incorporation of agentic AI within web browsers is altering information-seeking behaviors [14, 20, 25]; however, this development is solidifying a centralized, cloud-based infrastructure. Such an operational model depends on transmitting a continuous stream of sensitive user data—including browsing history, queries, and interactions—for external processing [21, 24]. Consequently, users

© ACM, 2026. This is the author’s version of the work.

The definitive version was published in: *Proceedings of the 2026 ACM SIGIR Conference on Human Information Interaction and Retrieval (CHIIR ’26)*, March 22–26, 2026, Seattle, WA, USA.

DOI: <https://doi.org/10.1145/3786304.3787913>

<sup>1</sup><https://github.com/saberzerhoudi/agentic-search-plugin>

are forced into a difficult position, having to select either sophisticated AI assistance or the control of their own data. In addition, using massive, general-purpose Large Language Models (LLMs) represents an inefficient use of computational resources for the specialized tasks in search assistance [6, 23]. This creates two central problems: first, a lack of privacy and user control, and second, an incorrect match between massive, high-cost models and the precise, contextual functions they are meant to perform.

In this work, we contend that recent technological developments offer a solution to this conflict. This solution is supported by two developments: first, the emergence of efficient, specialized Small Language Models (SLMs) suited for specific domains [1, 3, 6], and second, the new-found capability to run these SLMs directly in the browser using technologies such as WebGPU [9, 27]. By leveraging these two trends, we introduce a new framework that provides adaptive search assistance *while ensuring no data exits the user’s local device*. Our proposed solution is a hybrid, fully client-side architecture. Within this structure, a predictive probabilistic model learns the user’s policy, and its outputs are then used as a grounding mechanism for an SLM reasoner that operates within the browser [27].

This paper provides the following contributions:

- We design and implement a novel, fully in-browser architecture that grounds a local SLM reasoner with an adaptive, predictive probabilistic model.
- We introduce an in-browser online learning mechanism [8, 18] that effectively personalizes a generic user model to an individual’s specific search patterns using direct feedback.
- We demonstrate through a 3-week longitudinal study that this fully client-side framework measurably improves search efficiency without compromising user privacy or agency.
- We release our system as an open-source tool to facilitate further research in privacy-centric, in-browser agentic AI.

## 2 Related Work

Our work is positioned at the intersection of probabilistic user modeling, agent architectures, and in-browser language models.

### 2.1 Probabilistic User Models

Modeling search behavior using methods like click models and Markov models has proven effective for capturing aggregate user patterns [4, 10, 28, 31]. However, these models are traditionally predictive, not generative, and often lack mechanisms for real-time adaptation to an individual’s intent [11]. Our work extends this approach: we use a lightweight, adaptive probabilistic model as the *grounding mechanism* for a generative reasoning engine [32].

### 2.2 LLM-Driven Information Seeking

The application of LLMs to search tasks is promising [5, 32], but current methods rely on cloud-based APIs. This architecture reintroduces privacy risks [22, 24], and can produce ungrounded outputs detached from the user’s immediate context [15]. Our approach is distinct in that it *rejects* this cloud-dependent model by bringing the reasoning engine to the client, solving the privacy, and grounding problems simultaneously.

### 2.3 In-Browser and Small Language Models

Our work is made possible by two concurrent trends. First, the shift toward highly efficient, specialized Small Language Models (SLMs) [23] as a sustainable and effective solution for vertical domains like search assistance [1]. Second, advancements in web technologies (e.g., WebGPU [9, 16], WebLLM [27]) now permit these multi-billion parameter SLMs to be executed directly within the browser. While some in-browser models exist for recommendation [29, 33], they typically lack generative reasoning. To our knowledge, no prior work has unified an adaptive behavioral model with a local SLM reasoner in a fully client-side architecture for search assistance. Our framework is designed to fill this gap.

## 3 System Architecture and Methodology

We designed our framework to solve the challenge of creating a grounded, adaptive, and fully private search assistant. The entire architecture is implemented as a Firefox browser extension. The fundamental design principle of this system is to perform all computation—from data-logging to probabilistic modeling and generative inference—entirely on the client side. This method ensures that no sensitive user data leaves the user’s machine, thereby directly resolving the conflict between personalization and privacy.

The architecture combines the predictive accuracy of a lightweight probabilistic model with the generative reasoning capabilities of an in-browser Small Language Model (SLM). It operates as a passive assistant, observing user behavior and producing contextual suggestions. These suggestions are based on a dynamic, locally-stored understanding of the user’s state and evolving goals. The system is composed of three primary modules that work together:

- Behavioral Observation and Provenance.
- Dynamic User Modeling and State Estimation.
- Cognitive Inference and Suggestion.

### 3.1 Behavioral Observation and Provenance

The framework’s foundation is a perceptual layer that captures a user’s interaction data within an active search session. This module operates as a client-side listener, logging a curated set of high-level interaction events indicative of an information-seeking strategy:

- **Issued Queries:** The raw text of each query.
- **SERP Clicks:** The rank and URL of each click on a Search Engine Result Page (SERP).
- **Document Dwell Times:** The time spent on a document before returning to the SERP.

This data is stored locally using the browser’s IndexedDB API [12], a persistent, in-browser database. This design choice is central to our privacy-preserving commitment. To ensure complete user agency, the extension’s interface provides clear controls for the user to enable, disable, or clear all logged data. The user retains sole authority over their data.

### 3.2 Dynamic User Model

This module is the analytical core of the framework. It transforms raw interaction logs into an adaptive policy model that can predict the user’s subsequent actions. We formalize the search session as

a **Markov Decision Process (MDP)**, which provides a robust mathematical framework for modeling sequential decision-making.

*States and Actions.* A session is defined by states  $s_t \in S$  and actions  $a_t \in A$ . The state  $s_t$  represents the user’s current context (e.g., ViewingSERP, ReadingDocument). The action space  $A$  is discrete and includes strategic actions such as ClickDocument\_i (for  $i \in [1, 10]$ ), ReturnToSERP, and SubmitNewQuery(*type*), where *type* is classified as *generalization*, *specialization*, or *reformulation*.

*Model Initialization (Cold Start).* To address the cold-start problem and ensure the system provides immediate value, it is provisioned with pre-trained behavioral policies,  $\pi_\theta(a|s)$ . These policies function as representations of established searcher archetypes, such as “exploratory” or “lookup” [2]. We generated these initial policies by training a compact **Multi-Layer Perceptron (MLP)** [17] on clustered sessions from the AOL query log dataset [26]. The behavior-aware MLP was selected as the model architecture due to its minimal resource footprint and rapid inference capabilities, which makes it highly suitable for execution within the browser.

*Prediction.* The module’s primary function is to use the current policy  $\pi_\theta$  to compute a probability distribution over the next set of actions given the user’s current state  $s_t$ . The action with the highest probability,  $a^* = \arg \max_{a \in A} \pi_\theta(a|s_t)$ , is identified as the most probable next action. This prediction is the critical output used to ground the SLM.

A crucial component of this module is its capacity for in-browser online learning via direct user feedback. This mechanism allows the model to evolve from a generic archetype into a policy that is deeply personalized. When a user accepts ( $f_t = +1$ ) or discards ( $f_t = -1$ ) a suggestion, the policy network’s parameters  $\theta$  are updated using a simple and efficient policy gradient rule:

$$\theta_{t+1} = \theta_t + \alpha \cdot f_t \cdot \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) \quad (1)$$

Here,  $\alpha$  is the learning rate (set to 0.01). This update follows the REINFORCE algorithm [30], directly adjusting the likelihood of the action that produced the suggestion by either strengthening or weakening it, depending on the user’s explicit feedback. The entire process, from prediction to update, runs within a dedicated **Web Worker** [13]. This ensures that the model’s computational tasks do not block the browser’s main thread, thereby preventing any degradation of the user experience. The resulting updated policy parameters,  $\theta_{t+1}$ , are then saved locally in IndexedDB, which establishes a virtuous cycle of continuous personalization.

### 3.3 Hybrid Cognitive Inference

This module connects the probabilistic model’s abstract prediction to the generative, human-readable advice the user receives. The predicted next action,  $a^*$ , is *not* shown directly to the user. Instead, it serves as a *grounding directive* for the in-browser SLM.

This two-step process is a key design choice. Instead of prompting an unconstrained SLM with a raw history to “help the user”, our method first utilizes the probabilistic model to identify a high-likelihood *strategic* action,  $a^*$ . It is this action,  $a^*$ , that then guides the SLM’s reasoning process. This technique constrains the solution space, ensuring the generated advice is directly relevant to the user’s predicted behavior.

The predicted action is inserted into a prompt template along with recent session history. A simplified prompt structure is:

You are a helpful search assistant. Based on the user’s recent activity, they will likely perform the following action: **{Predicted Action  $a^*$ }**.  
The user’s last query was “**{Last Query}**”.  
Generate a concise, helpful suggestion to guide them.

The entire inference pipeline runs locally. We leverage the **WebGPU engine** to execute a ~2.7B parameter model from the Phi family [1] directly in the browser. This model was selected as it provides state-of-the-art balanced high performance with a suitable footprint for client-side execution [19]. Our tests confirmed the viability of this approach on modern hardware<sup>2</sup>. The local SLM generated suggestions with an average latency of ~9 tokens/sec, which is well within the threshold for a non-disruptive user experience.

Finally, to avoid “suggestion fatigue”, the system incorporates an *intrusion avoidance* heuristic. It refrains from generating suggestions when (1) the model predicts a clear navigational action with high confidence, or (2) the probability distribution over next actions is nearly uniform, indicating high model uncertainty.

## 4 Experimental Evaluation

Our evaluation was designed as a longitudinal study to validate the framework’s core claims. The goals were to (1) quantitatively measure the effectiveness of the in-browser personalization mechanism and (2) assess the framework’s real-world impact on both the search behavior and subjective experience of users.

### 4.1 Experimental Setup

*Participants.* We recruited 18 participants (12 male, 6 female, ages 21-32) from a local university, representing a range of technical backgrounds with high self-reported familiarity with web search.

*Procedure.* We conducted a 3-week longitudinal study, which was critical for observing the model’s adaptation over time.

- **Week 1 (Baseline):** Participants installed the extension with suggestions *disabled*. This allowed us to collect baseline data on their natural, un-assisted search behavior.
- **Weeks 2-3 (Intervention):** The suggestion feature was *enabled*. Participants continued their normal web tasks, interacting with the suggestions as they chose.

*Data Collection.* To ensure a realistic evaluation environment, participants used their own laptops running Firefox. The extension logged all anonymized interaction data locally. At the study’s conclusion, participants *themselves* exported and submitted this anonymized data. They also completed a post-study questionnaire assessing usability (SUS) [7], perceived utility, and trust.

### 4.2 Model Personalization and Adaptation

*Motivation.* Our first research question was to determine if the in-browser online learning mechanism could effectively adapt a

<sup>2</sup>Apple MacBook Pro (M1 Pro, 10-core CPU / 16-core GPU, 32 GB unified memory).

**Table 1: Next-action prediction performance. The personalized in-browser model outperforms the static, generic baselines.**

Model	Pred. Acc. (%)	MRR
Generic-Exploratory	31.2	0.39
Generic-Lookup	26.5	0.34
<b>Persona-Adapted</b>	<b>38.7</b>	<b>0.47</b>

generic model to an individual’s unique search patterns. We hypothesized that a personalized model would outperform the initial generic models in predicting that user’s *own* future actions.

*Setting.* We analyzed the interaction logs from the intervention phase using a temporal hold-out strategy: for each participant, the first 80% of their interactions were used to simulate the adaptation process (training via Eq. 1), and the final 20% served as the test set for a next-action prediction task.

*Baselines.* We compared the predictive performance of the pre-trained generic models (Generic-Exploratory, Generic-Lookup) against the personalized **Persona-Adapted** model for each user.

*Results.* The results, summarized in Table 1, provide strong quantitative support for our hypothesis. The Persona-Adapted model (38.7% Accuracy, 0.47 MRR) outperformed both generic baselines. Compared to the strongest baseline (Generic-Exploratory), our adapted model achieved a 24.0% relative improvement in prediction accuracy and 20.5% in MRR. This confirms that the online learning mechanism is successfully capturing user-specific search habits, evolving from a generic archetype into a specialized policy.

### 4.3 Behavioral Shift and User Perception

*Motivation.* Our ultimate goal is not just to predict actions, but to *improve* the user’s search process. We evaluated whether exposure to the assistant led to measurable changes in search behavior and assessed the user’s subjective experience.

*Setting.* We compared behavioral metrics from the one-week Baseline phase (no suggestions) against the two-week Intervention phase using paired *t*-tests. We also analyzed the subjective ratings from the post-study questionnaire.

*Results.* The analysis in Table 2 revealed a statistically positive improvement in user search behavior. Participants’ queries became measurably more complex during the intervention (4.1 vs. 3.5 terms,  $p < 0.05$ ), suggesting the assistant guided them to more descriptive formulations. This, in turn, led to a decrease in average session length (5.2 vs. 6.8 queries,  $p < 0.05$ ), indicating a clear gain in search efficiency.

The mean suggestion acceptance rate was 36.4%. This result is notable as it shows users maintained their agency, evaluating suggestions rather than passively accepting them. Subjective feedback was highly positive, with an “Excellent” SUS score of 82.5. Critically for a privacy-focused tool, users reported both high perceived utility (6.1/7) and trust (5.9/7).

**Table 2: Comparison of behavioral and subjective metrics between the Baseline and Intervention phases. Asterisks (\*) denote a statistically significant difference ( $p < 0.05$ ).**

Metric	Baseline	Intervention
Avg. Query Complexity (terms)	3.5	4.1*
Avg. Session Length (queries)	6.8	5.2*
Suggestion Acceptance Rate (%)	-	36.4
System Usability Scale (SUS)	-	82.5
Perceived Utility (1–7)	-	6.1
Trust in Suggestions (1–7)	-	5.9

## 5 Discussion and Conclusion

Our experimental evaluation confirms the framework’s effectiveness. The in-browser online learning mechanism successfully adapted the generic policy, yielding a 24.0% improvement in next-action prediction accuracy over the strongest baseline. This personalization also measurably improved search efficiency: we observed a statistically reduction in average session length (from 6.8 to 5.2 queries) as users were guided toward more effective, complex search strategies.

A 36.4% suggestion acceptance rate further validates our core architectural thesis. It demonstrates that our hybrid approach—combining a probabilistic model’s behavioral predictions with an in-browser SLM’s reasoning—provides useful, context-aware guidance while preserving complete user control.

Most importantly, this work challenges the prevailing assumption that sophisticated agentic assistance must be centralized. Our findings demonstrate a viable, computationally efficient, and fully private alternative. We show that by using specialized SLMs with in-browser execution, it is possible to achieve effective personalization without compromising user data sovereignty.

### 5.1 Limitations and Future Work

We acknowledge this study’s limitations. The participant sample was drawn from a university population and is not representative of the general public. Furthermore, our evaluation was designed to isolate the effect of personalization, comparing the adapted model only against its own generic baseline rather than external methods.

Future work will focus on benchmarking our adaptive model against more complex sequential recommendation baselines to situate its predictive performance. We also plan to extend the framework to process multimodal content and to integrate more advanced reinforcement learning algorithms for policy adaptation, all while maintaining the fully client-side, privacy-preserving architecture.

### 5.2 Conclusion

This paper presented a novel framework for a privacy-preserving search assistant that operates entirely on the user’s device. We designed a system that learns from and adapts to the user by grounding an in-browser SLM with an adaptive policy model. Our results show a viable method for providing sophisticated AI assistance while maintaining complete user autonomy. We provide our framework as an open-source tool to support further research and development in this user-centric, in-browser paradigm.

## References

- [1] Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Björck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, and Ronen Eldan et al. 2024. Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone. arXiv:2404.14219 [cs.CL] <https://arxiv.org/abs/2404.14219>
- [2] Manoj K. Agarwal and Tezan Sahu. 2021. Lookup or Exploratory: What is Your Search Intent? arXiv:2110.04640 [cs.IR] <https://arxiv.org/abs/2110.04640>
- [3] Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Becked, Lewis Tunstall, Andrés Marafioti, Hynek Kydlíček, Agustín Piqueres Lajarin, Vaibhav Srivastav, Joshua Lochner, Caleb Fahlgrén, Xuan-Son Nguyen, Clémentine Fourrier, Ben Burtenshaw, Hugo Larcher, Haojun Zhao, Cyril Zakka, Mathieu Morlon, Colin Raffel, Leandro von Werra, and Thomas Wolf. 2025. SmolLM2: When Smol Goes Big – Data-Centric Training of a Small Language Model. arXiv:2502.02737 [cs.CL] <https://arxiv.org/abs/2502.02737>
- [4] Mamoun A. Awad and Issa Khalil. 2012. Prediction of user’s web-browsing behavior: Application of markov model. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* 42, 4 (2012), 1131–1142.
- [5] Jonas Becker. 2024. Multi-Agent Large Language Models for Conversational Task-Solving. arXiv:2410.22932 [cs.CL] <https://arxiv.org/abs/2410.22932>
- [6] Peter Belcak, Greg Heinrich, Shizhe Diao, Yonggan Fu, Xin Dong, Saurav Muralidharan, Yingyan Celine Lin, and Pavlo Molchanov. 2025. Small Language Models are the Future of Agentic AI. arXiv:2506.02153 [cs.AI] <https://arxiv.org/abs/2506.02153>
- [7] John Brooke. 1996. SUS: A “Quick and Dirty” Usability Scale. In *Usability Evaluation in Industry*, P. W. Jordan, B. Thomas, I. L. McClelland, and B. Weerdmeester (Eds.). Taylor & Francis, 189–194.
- [8] Davide Cacciarrelli and Murat Kulahci. 2023. Active learning for data streams: a survey. *Machine Learning* 113, 1 (Nov. 2023), 185–239. doi:10.1007/s10994-023-06454-2
- [9] Zhiyang Chen, Yun Ma, Haiyang Shen, and Mugeng Liu. 2025. WeInfer: Unleashing the Power of WebGPU on LLM Inference in Web Browsers. In *Proceedings of the ACM on Web Conference 2025, WWW 2025, Sydney, NSW, Australia, 28 April 2025 - 2 May 2025*, Guodong Long, Michale Blumstein, Yi Chang, Liane Lewin-Eytan, Zi Helen Huang, and Elad Yom-Tov (Eds.). ACM, 4264–4273. <https://doi.org/10.1145/3696410.3714553>
- [10] Flavio Chierichetti, Ravi Kumar, Prabhakar Raghavan, and Tamás Szilárd. 2012. Are web users really Markovian?. In *Proceedings of the 21st World Wide Web Conference 2012, WWW 2012, Lyon, France, April 16-20, 2012*, Alain Mille, Fabien Gandon, Jacques Misselis, Michael Rabinovich, and Steffen Staab (Eds.). ACM, 609–618. doi:10.1145/2187836.2187919
- [11] Siamak Farshidi, Kiyan Rezaee, Sara Mazaheri, Amir Hossein Rahimi, Ali Dadashzadeh, Morteza Ziabakhsh, Sadegh Eskandari, and Slinger Jansen. 2023. Understanding User Intent Modeling for Conversational Recommender Systems: A Systematic Literature Review. arXiv:2308.08496 [cs.IR] <https://arxiv.org/abs/2308.08496>
- [12] Web Applications Working Group. 2015. *Indexed Database API*. W3C Recommendation REC-IndexedDB-20150108. World Wide Web Consortium (W3C). <https://www.w3.org/TR/2015/REC-IndexedDB-20150108/>
- [13] Web Applications Working Group. 2021. *Web Workers*. W3C Working Group Note NOTE-workers-20210128. World Wide Web Consortium (W3C). <https://www.w3.org/TR/2021/NOTE-workers-20210128/>
- [14] Google Inc. 2024. *Generative AI in Search: Let Google do the searching for you*. <https://blog.google/products/search/generative-ai-google-search-may-2024/> Google Blog.
- [15] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of Hallucination in Natural Language Generation. *Comput. Surveys* 55, 12 (March 2023), 1–38. doi:10.1145/3571730
- [16] Benjamin Kenwright. 2022. Introduction to the webgpu api. In *Acm siggraph 2022 courses*. 1–184.
- [17] Weixin Li, Yuhao Wu, Yang Liu, Weike Pan, and Zhong Ming. 2024. BMLP: Behavior-aware MLP for Heterogeneous Sequential Recommendation. arXiv:2402.12733 [cs.IR] <https://arxiv.org/abs/2402.12733>
- [18] Yuxuan Lu, Jing Huang, Hui Liu, Jiri Gesi, Yan Han, Shiha Fu, Tianqi Zheng, and Dakuo Wang. 2025. WEBSERV: A Browser-Server Environment for Efficient Training of Reinforcement Learning-based Web Agents at Scale. arXiv:2510.16252 [cs.LG] <https://arxiv.org/abs/2510.16252>
- [19] Zhenyan Lu, Xiang Li, Dongqi Cai, Rongjie Yi, Fangming Liu, Wei Liu, Jian Luan, Xiwen Zhang, Nicholas D. Lane, and Mengwei Xu. 2025. Demystifying Small Language Models for Edge Deployment. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (Eds.). Association for Computational Linguistics, 14747–14764. <https://aclanthology.org/2025.acl-long.718/>
- [20] Yusuf Mehdi. 2023. *Reinventing search with a new AI-powered Microsoft Bing and Edge, Your Copilot for the Web*. <https://blogs.microsoft.com/blog/2023/02/07/reinventing-search-with-a-new-ai-powered-microsoft-bing-and-edge-your-copilot-for-the-web/> Microsoft Corporation Blog.
- [21] Arvind Narayanan and Vitaly Shmatikov. 2008. Robust De-anonymization of Large Sparse Datasets. In *2008 IEEE Symposium on Security and Privacy (SP 2008), 18-21 May 2008, Oakland, California, USA*. IEEE Computer Society, 111–125. doi:10.1109/SP.2008.33
- [22] Arvind Narayanan and Vitaly Shmatikov. 2008. Robust De-anonymization of Large Sparse Datasets. In *2008 IEEE Symposium on Security and Privacy (SP 2008), 18-21 May 2008, Oakland, California, USA*. IEEE Computer Society, 111–125. doi:10.1109/SP.2008.33
- [23] Chien Van Nguyen, Xuan Shen, Ryan Aponte, Yu Xia, Samyadeep Basu, Zhengmian Hu, Jian Chen, Mihir Parmar, Sasidhar Kunapuli, Joe Barrow, Junda Wu, Ashish Singh, Yu Wang, Jiuxiang Gu, Franck Dernoncourt, Nesreen K. Ahmed, Nedim Lipka, Ruiyi Zhang, Xiang Chen, Tong Yu, Sungchul Kim, Hanieh Deilamsalehy, Namyong Park, Mike Rimer, Zhehao Zhang, Huanrui Yang, Ryan A. Rossi, and Thien Huu Nguyen. 2024. A Survey of Small Language Models. arXiv:2410.20011 [cs.CL] <https://arxiv.org/abs/2410.20011>
- [24] Paul Ohm. 2009. Broken promises of privacy: Responding to the surprising failure of anonymization. *UCLA L. Rev.* 57 (2009), 1701.
- [25] OpenAI. 2025. *Introducing ChatGPT Atlas*. <https://openai.com/index/introducing-chatgpt-atlas/> OpenAI Blog.
- [26] Gilad Pass, Abdur Chowdhury, and Chris Torgeson. 2006. AOL Query Log: 20 Million Queries from 650 000 Users (March–May 2006). Web research dataset released by AOL Research. <https://www.kaggle.com/datasets/dineshdyv/aol-user-session-collection-500k> Distributed as “AOL User Session Collection 500K”; anonymised; for research use only.
- [27] Charlie F. Ruan, Yucheng Qin, Xun Zhou, Ruihang Lai, Hongyi Jin, Yixin Dong, Bohan Hou, Meng-Shiun Yu, Yiyan Zhai, Sudeep Agarwal, Hangrui Cao, Siyuan Feng, and Tianqi Chen. 2024. WebLLM: A High-Performance In-Browser LLM Inference Engine. arXiv:2412.15803 [cs.LG] <https://arxiv.org/abs/2412.15803>
- [28] Vu T Tran and Norbert Fuhr. 2013. Markov modeling for user interaction in retrieval. In *SIGIR 2013 Workshop on Modeling User Behavior for Information Retrieval Evaluation (MUBE 2013)*, Vol. 5. Citeseer, 7–2.
- [29] Qipeng Wang, Shiqi Jiang, Zhenpeng Chen, Xu Cao, Yuanchun Li, Aoyu Li, Yun Ma, Ting Cao, and Xuanzhe Liu. 2024. Anatomizing Deep Learning Inference in Web Browsers. arXiv:2402.05981 [cs.LG] <https://arxiv.org/abs/2402.05981>
- [30] Ronald J. Williams. 1992. Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement Learning. *Machine Learning* 8, 3–4 (1992), 229–256. doi:10.1007/BF00992696
- [31] Saber Zerhouni, Michael Granitzer, Christin Seifert, and Jörg Schlöter. 2022. Simulating User Interaction and Search Behaviour in Digital Libraries. In *Proceedings of the 18th Italian Research Conference on Digital Libraries, Padua, Italy, February 24-25, 2022 (hybrid event) (CEUR Workshop Proceedings, Vol. 3160)*, Giorgio Maria Di Nunzio, Beatrice Portelli, Domenico Redavid, and Gianmaria Silvello (Eds.). CEUR-WS.org. <https://ceur-ws.org/Vol-3160/paper8.pdf>
- [32] Erhan Zhang, Xingzhu Wang, Peiyuan Gong, Yankai Lin, and Jiaxin Mao. 2024. USimAgent: Large Language Models for Simulating Search Users. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2024, Washington DC, USA, July 14-18, 2024*, Grace Hui Yang, Hongning Wang, Sam Han, Claudia Hauff, Guido Zuccon, and Yi Zhang (Eds.). ACM, 2687–2692. doi:10.1145/3626772.3657963
- [33] Wenchao Zhao, Xiaoyi Liu, Ruilin Xu, Lingxi Xiao, and Muqing Li. 2024. E-commerce Webpage Recommendation Scheme Based on Semantic Mining and Neural Networks. *Journal of Theory and Practice of Engineering Science* 4, 03 (March 2024), 207–215. doi:10.53469/jtpes.2024.04(03).20