# Modeling Spoken Information Queries for Virtual Assistants

## Open Problems, Challenges and Opportunities

Christophe Van Gysel
cvangysel@apple.com
Apple
Cambridge, MA, USA

## ABSTRACT

Virtual assistants are becoming increasingly important speech-driven Information Retrieval platforms that assist users with various tasks. We discuss open problems and challenges with respect to modeling spoken information queries for virtual assistants, and list opportunities where Information Retrieval methods and research can be applied to improve the quality of virtual assistant speech recognition. We discuss how query domain classification, knowledge graphs and user interaction data, and query personalization can be helpful to improve the accurate recognition of spoken information domain queries. Finally, we also provide a brief overview of current problems and challenges in speech recognition.

## CCS CONCEPTS

• **Information systems** → **Search interfaces**; **Query log analysis**; • **Computing methodologies** → **Speech recognition**.

## KEYWORDS

virtual assistants, query log analysis, automated speech recognition

## 1 INTRODUCTION

Virtual assistants (VAs) are becoming increasingly important [27] Information Retrieval (IR) platforms that assist users with various tasks. Users primarily interact with VAs through voice commands, where users initiate a retrieval request by uttering a query, possibly preceded by a wake word (e.g., *"hey VA"*). Accurately transcribing spoken voice queries [13], for subsequent processing by the retrieval engine, is a challenging problem that can benefit greatly from knowledge of the IR application.

Automated Speech Recognition (ASR) systems, responsible for transcribing the spoken utterance, are trained on audio/text pairs that are expensive to obtain. Language models (LMs) are a component within ASR systems that act as a query prior and are trained

on only text. The LM becomes increasingly important when the spoken query is ambiguous or difficult to understand (e.g., unintelligible speech). Take as an example the encyclopedia query *"what is borrelia"* where the user intends to obtain information about the Borrelia bacteria. In this particular case, the entity *borrelia* may be misrecognized as *gorilla* if the LM assigns a low likelihood to the conditional probability P ( *borrelia* | what is). However, this problem can be alleviated through query log analysis, injection of external knowledge (e.g., entity popularity), and use of contextual signals, amongst other methods.

The challenging nature of VA query recognition is further exacerbated by stringent runtime requirements. Recognition needs to occur in real-time as users expect results soon after they finish speaking. When ASR occurs on-device, model size becomes an additional constraint—since disk space and network bandwidth are costly. In addition, the ability to patch or perform incremental updates of models is desirable functionality. Finally, LMs need to be trained within a reasonable amount of time, since otherwise they may be outdated by the time the LMs reach edge devices.

We provide a succinct overview of the use of LMs in ASR, and subsequently cover topics on using knowledge of the IR application to improve ASR with a focus on entity-heavy VA queries: (1) query domain classification, (2) entity popularity and knowledge graph (KG) mining, and (3) personalization. Finally, we briefly cover non-IR topics relating to LMs and ASR.

## 2 ASR PRIMER

Automated Speech Recognition is the task of translating a speech signal $X$ into a string of words $s$. Contemporary ASR systems can be divided into two categories: (a) traditional **hybrid systems** that rely on Bayes' rule to combine acoustic and language model components, and (b) more modern **end-to-end systems** that directly predict output sequences of text from acoustic representations.

### 2.1 Hybrid ASR systems

Hybrid ASR systems operate by decomposing the ASR task using Bayes' rule as follows [18, p. 289]:

$$S^* = \text{argmax}_s P(s \mid X) = \text{argmax}_s \frac{P(X \mid s)}{P(X)} \cdot P(s)$$
$$= \text{argmax}_s P(X \mid s) \cdot P(s), \quad (1)$$

where $P(X \mid s)$ is provided by the acoustic model (AM) and denotes the likelihood of speech signal $X$ given the string of words $s$, and $P(s)$ is provided by the LM and denotes the prior probability of a string of words. $P(X)$ is the probability of the speech signal and can be ignored as it is constant for all hypotheses. The AM and LM are trained independently and subsequently combined.

## 2.2 End-to-end ASR systems

End-to-end ASR (E2E) systems directly compute the probability distribution $P(s \mid X)$ of output strings of words $s$ given speech signals $X$, and are typically implemented using neural encoder-decoder architectures [19, §16.3]. Since E2E systems model the ASR task directly, they are trained on paired audio-text data–which can be expensive to obtain and may not provide full coverage for tail utterances. Hence, often, an additional LM, trained only on abundantly available texts, and external to the E2E model, is combined through interpolation as follows:

$$S^* = \operatorname{argmax}_s P(s \mid X) \approx \operatorname{argmax}_s P_{\text{E2E}}(s \mid X) \cdot P_{\text{Ext. LM}}(s)^{\lambda}, \quad (2)$$

where $P_{\text{E2E}}(s \mid X)$ is provided by the E2E model, $P_{\text{Ext. LM}}(s)$ is the LM probability and $\lambda$ is an interpolation hyperparameter.

## 2.3 Language models

Regardless of ASR system architecture, hybrid or E2E, a LM trained solely on text data can be used to improve recognition quality. The LM builds on the chain rule of probability:

$$P(W) = P(w_1 w_2 \dots w_n) = \prod_{i=1}^{N} P(w_i \mid w_1 w_2 \dots w_{i-1}). \quad (3)$$

In practice, strings of words are wrapped in special start/end markers (<s>/</s>, resp.) to denote the beginning and the end of the string of words (which is, typically, a sentence). For example, the prior probability of utterance *SIGIR* would be computed as

$$P(\text{<s> SIGIR </s>}) = P(\text{SIGIR} \mid \text{<s>}) P(\text{</s>} \mid \text{<s> SIGIR}),$$

where <s> and </s> mark the beginning and end of sentence, resp.

## 3 OPEN PROBLEMS AND CHALLENGES

## 3.1 Use of query domain classifications

VA queries can be categorized according to domains where each domain supports specific use-cases. For example, there exist media player queries such as *"play the look by metronomy"* where the user instructs the VA to play the song "The Look" by the band Metronomy, or encyclopedic queries where the user wants to learn more about a specific entity (e.g., *"who is joe biden"*).

In this section we discuss the application of query domain classifications, possibly provided by NLP/IR methods, to improve VA ASR, either by (a) using domain classifications at runtime to guide the ASR decoding process, or (b) utilizing the classification of queries at LM training time.

*3.1.1 Improving the ASR decoding process.* At recognition time, contextual signals, such as partial recognition hypotheses [25] or the user location [39], can be used to modify the search space. Pusateri et al. [25] combine multiple domain-specific expert n-gram LMs into a single LM by weighing the expert LMs based on the confidence expressed by each expert LM on how well they support specific left spoken contexts. Following the example above, the media player domain LM would receive a large weight following the left context *"<s> play"*, whereas a LM trained on encyclopedia queries may be well-suited following left context *"<s> who is"*.
**Relevance to IR.** From an IR perspective, contextual signals extracted from (partial) user interactions (e.g., session information,

partial queries) with the VA can be integrated into the ASR component responsible for combining multiple domain-specific expert models. Effective integration of contextual signals into E2E ASR systems (i.e., Eq. 2) remains an open problem today.

*3.1.2 Building better LMs by leveraging query domain classifications.* Gondala et al. [10] take a different approach and use classifications of training data queries to influence the n-gram LM training algorithm. For example, query domains that reference many tail entities can be allocated more model capacity and that in turn improves the recognition of tail entities.
**Relevance to IR.** Offline classification of query logs can be used to improve ASR. While in [10], the authors used a domain-driven generative process [35] to obtain training query texts, their approach can also be applied on queries that occur in usage logs. However, ASR is a noisy process and consequently queries may contain recognition errors, and hence, domain classification methods designed for typed query traffic may not directly apply. The IR community may find the usage of signals made available during the ASR decoding process, such as word-level confidence [17], helpful to adapt methods designed for typed query classification to spoken queries.

## 3.2 KGs and other external data

As mentioned at the end of the previous section, spoken query logs contain recognition errors, and LMs used for ASR are often trained, at least partially, on query logs. This practice can lead to a feedback loop of reinforced errors. While filtering techniques can provide some relief, they may also introduce undesirable biases in the training data. The use of external data sources is an alternative solution that can benefit the recognition of entity-rich queries.

*3.2.1 Query templatization and entity popularity.* Gandhe et al. [9] estimate n-gram LMs directly from entity-rich grammars to improve ASR for new application intents in VAs. In this case, queries such as in the example of §3.1 can be represented as templates (e.g., *"play $SONG by $ARTIST"*) with entity slots. Van Gysel et al. [35] released a VA media player query grammar, including a large list of media player entities extracted from a large-scale media catalog user interactions. In [36], the authors extract entities from a VA query log that occur in the presence of spoken left context (e.g., a verb) to improve the recognition of entity name queries [42] in the absence of left context.
**Relevance to IR.** From an IR point of view, there exist multiple challenges. First of all, while query templates can be created manually by domain experts, methods to automatically extract templates from a query log can be useful. Secondly, while entity popularity can be extracted from external sources, there likely still exists a gap between popularity in the source application and the VA application (e.g., a difference in demographics). Hence, entity popularity adaptation methods are still an open research problem.

*3.2.2 Using KG relations during ASR decoding.* Saebi et al. [28], amongst others [15, 22], make use of entity type and entity–entity relations during the ASR decoding process to improve the recognition of tail named entities. For example, if the ASR decoder is considering two hypotheses (a) *"play can you moon by Harry Styles"* and (b) *"play Canyon Moon by Harry Styles"*, their approach will

use the KG relationship between artist (i.e., "Harry Styles") and song title (i.e., "Canyon Moon") as a signal during recognition to boost the likelihood that the factually correct hypothesis (b) is chosen.

**Relevance to IR.** Hence, IR research focusing on improving KGs and entity linking in spoken queries can directly improve the effectiveness of VA ASR.

### 3.3 Personalization

Personalization of on-device ASR is an active area of research [3]. For the VA application, and from the language modeling perspective (as opposed to acoustics [34]), systems may be able to benefit from signals used in other search applications, such as Web search [32], as users with different profiles tend to search for different sets of topics. More specifically, knowledge about the user's interests–which may eventually lead to an interaction with a specific intent–can be helpful to improve user experience. Xiao et al. [39] improve ASR by bucketing users according to their coarse geographic location and enable region-specific query LMs during the ASR decoding process. By personalizing the ASR query model based on user location, they show a significant improvement in the accurate recognition of spoken point-of-interest queries.

**Relevance to IR.** On the IR side, user models [7] based on query behavior or other signals may be helpful to power futher personalization of on-device VA ASR.

### 3.4 Beyond IR

In the previous sections, we focused on the impact of IR research on the accurate recognition of spoken information queries for the VA application. Naturally, there exist a multitude of challenges on the ASR side as well. End-to-end ASR models [38], as opposed to traditional Gaussian mixture models, have been increasingly gaining popularity since end-to-end models consist of less components—hence, reducing maintenance costs. However, integration of external LMs into [5, 21, 29], and personalization of [11, 33, 34], end-to-end systems remains an active research area. With respect to LM, Neural Network LMs (NNLM) [1] have gained popularity within ASR [12, 30, 43]. In the case of VA ASR, NNLMs can be significantly more economical in terms of storage costs than their N-Gram LM [20] counterparts, as with the latter, the size of the models grows proportional to the data complexity. However, practical limitations such as training an NNLM from a heterogenous corpora, inference latency [26], and federated learning [41] remain challenging. More recently, large pre-trained Transformer LMs [4, 8, 23, 37] have also been used to improve ASR [6, 16, 31, 40], although a domain gap may exist [24]; and have also been shown to be effective for synthetic data generation [2, 14].

## 4 CONCLUSIONS

We discussed open problems and challenges with respect to modeling spoken information queries for VAs, and listed opportunities where IR methods and research can be applied to improve the quality of VA ASR. More specifically, we discussed how query domain classification can be used during speech recognition and to build better LMs. Next, we discussed the use of KGs and external data sources based on user interactions, and discussed personalization.

Finally, and for completeness, we briefly provided an overview of challenges and open problems within ASR.

We hope that the discussed topics are useful to IR researchers and lead to the exploration of new, cross-disciplinary research directions and ideas, and as inspiration to discover new application domains for existing methods.

## SPEAKER BIOGRAPHY

Christophe Van Gysel is a Staff Research Scientist working on the Siri Speech language modeling team at Apple where he works on the boundary between ASR and Search. Christophe obtained his PhD in Computer Science from the University of Amsterdam in 2017. During his PhD, Christophe worked on neural ranking using representation learning models with a focus on entities and published at WWW, SIGIR, CIKM, WSDM, TOIS, amonst others.

## COMPANY PROFILE

Apple revolutionised personal technology with the introduction of the Macintosh in 1984. Today, Apple leads the world in innovation with iPhone, iPad, Mac, Apple Watch, and Apple TV. Apple's five software platforms — iOS, iPadOS, macOS, watchOS, and tvOS — provide seamless experiences across all Apple devices and empower people with breakthrough services including the App Store, Apple Music, Apple Pay, and iCloud. Apple's more than 100,000 employees are dedicated to making the best products on earth, and to leaving the world better than we found it.

## REFERENCES

[1] Yoshua Bengio, Réjean Ducharme, and Pascal Vincent. 2000. A neural probabilistic language model. *NeurIPS* (2000).

[2] Luiz Bonifacio, Hugo Abonizio, Marzieh Fadaee, and Rodrigo Nogueira. 2022. Inpars: Unsupervised dataset generation for information retrieval. In *SIGIR*. 2387–2392.

[3] Theresa Breiner, Swaroop Ramaswamy, Ehsan Variani, Shefali Garg, Rajiv Mathews, Khe Chai Sim, Kilol Gupta, Mingqing Chen, and Lara McConnaughey. 2022. UserLibri: A Dataset for ASR Personalization Using Only Text. In *Interspeech*.

[4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in Neural Information Processing Systems* 33 (2020), 1877–1901.

[5] Rodrigo Cabrera, Xiaofeng Liu, Mohammadreza Ghodsi, Zebulun Matteson, Eugene Weinstein, and Anjuli Kannan. 2021. Language model fusion for streaming end to end speech recognition. *arXiv preprint arXiv:2104.04487* (2021).

[6] Shih-Hsuan Chiu and Berlin Chen. 2021. Innovative BERT-based reranking language models for speech recognition. In *SLT*. IEEE, 266–271.

[7] Aleksandr Chuklin, Ilya Markov, and Maarten de Rijke. 2015. Click models for web search. *Synthesis lectures on information concepts, retrieval, and services* 7, 3 (2015), 1–115.

[8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL-HTL*. 4171–4186.

[9] Ankur Gandhe, Ariya Rastrow, and Bjorn Hoffmeister. 2018. Scalable Language Model Adaptation for Spoken Dialogue Systems. In *SLT*. IEEE, 907–912.

[10] Sashank Gondala, Lyan Verwimp, Ernest Pusateri, Manos Tsagkias, and Christophe Van Gysel. 2021. Error-Driven Pruning of Language Models for Virtual Assistants. In *ICASSP*. IEEE, 7413–7417.

[11] Aditya Gourav, Linda Liu, Ankur Gandhe, Yile Gu, Guitang Lan, Xiangyang Huang, Shashank Kalmane, Gautam Tiwari, Denis Filimonov, Ariya Rastrow, et al. 2021. Personalization strategies for end-to-end speech recognition systems. In *ICASSP*. IEEE, 7348–7352.

[12] Alex Graves, Navdeep Jaitly, and Abdel-rahman Mohamed. 2013. Hybrid speech recognition with deep bidirectional LSTM. In *ASRU*. IEEE, 273–278.

[13] Ido Guy. 2016. Searching by talking: Analysis of voice queries on mobile web search. In *SIGIR*. 35–44.

[14] Perttu Hämäläinen, Mikke Tavast, and Anton Kunnari. 2023. Evaluating Large Language Models in Generating Synthetic HCI Research Data: a Case Study. In *CHI*.

[15] Hiroaki Hayashi, Zecong Hu, Chenyan Xiong, and Graham Neubig. 2020. Latent relation language models. In *AAAI*. 7911–7918.

[16] Hongzhao Huang and Fuchun Peng. 2019. An empirical study of efficient ASR rescoring with transformers. *arXiv preprint arXiv:1910.11450* (2019).

[17] Woojay Jeon, Maxwell Jordan, and Mahesh Krishnamoorthy. 2020. On Modeling ASR Word Confidence. In *ICASSP*. IEEE, 6324–6328.

[18] Daniel Jurafsky and James H. Martin. 2008. *Speech and Language Processing, 2nd edition.* Prentice Hall.

[19] Daniel Jurafsky and James H. Martin. 2023. *Speech and Language Processing, 3rd edition (January 7, 2023 draft).* Prentice Hall.

[20] Slava Katz. 1987. Estimation of probabilities from sparse data for the language model component of a speech recognizer. *IEEE transactions on acoustics, speech, and signal processing* (1987), 400–401.

[21] Suyoun Kim, Yuan Shangguan, Jay Mahadeokar, Antoine Bruguier, Christian Fuegen, Michael L. Seltzer, and Duc Le. 2021. Improved Neural Language Model Fusion for Streaming Recurrent Neural Network Transducer. In *ICASSP*. 7333–7337.

[22] Robert L Logan IV, Nelson F Liu, Matthew E Peters, Matt Gardner, and Sameer Singh. 2019. Barack's wife hillary: Using knowledge-graphs for fact-aware language modeling. In *ACL*.

[23] OpenAI. 2023. *GPT-4 Technical Report.* Technical Report. OpenAI.

[24] Valentin Pelloin, Franck Dary, Nicolas Hervé, Benoît Favre, Nathalie Camelin, Antoine Laurent, and Laurent Besacier. 2022. ASR-Generated Text for Language Model Pre-training Applied to Speech Tasks. *arXiv preprint arXiv:2207.01893* (2022).

[25] Ernest Pusateri, Christophe Van Gysel, Rami Botros, Sameer Badaskar, Mirko Hannemann, Youssef Oualil, and Ilya Oparin. 2019. Connecting and Comparing Language Model Interpolation Techniques. In *Interspeech*. 3500–3504.

[26] Anirudh Raju, Denis Filimonov, Gautam Tiwari, Guitang Lan, and Ariya Rastrow. 2019. Scalable multi corpora neural language models for asr. In *Interspeech*.

[27] Juniper Research. 2019. Digital Voice Assistants in Use to Triple to 8 Billion by 2023, Driven by Smart Home Devices. Press Release.

[28] Mandana Saebi, Ernest Pusateri, Aaksha Meghawat, and Christophe Van Gysel. 2021. A Discriminative Entity-Aware Language Model for Virtual Assistants. In *Interspeech*.

[29] Changhao Shan, Chao Weng, Guangsen Wang, Dan Su, Min Luo, Dong Yu, and Lei Xie. 2019. Component fusion: Learning replaceable language model component for end-to-end speech recognition system. In *ICASSP*. IEEE, 5361–5635.

[30] Apeksha Shewalkar, Deepika Nyavanandi, and Simone A Ludwig. 2019. Performance evaluation of deep neural networks applied to speech recognition: RNN, LSTM and GRU. *Journal of Artificial Intelligence and Soft Computing Research* 9, 4 (2019), 235–245.

[31] Joonbo Shin, Yoonhyung Lee, and Kyomin Jung. 2019. Effective sentence scoring method using bert for speech recognition. In *Asian Conference on Machine Learning*. 1081–1093.

[32] Ahu Sieg, Bamshad Mobasher, and Robin Burke. 2007. Web search personalization with ontological user profiles. In *Proceedings of the sixteenth ACM conference on Conference on information and knowledge management*. 525–534.

[33] Khe Chai Sim, Françoise Beaufays, Arnaud Benard, Dhruv Guliani, Andreas Kabel, Nikhil Khare, Tamar Lucassen, Petr Zadrazil, Harry Zhang, Leif Johnson, et al. 2019. Personalization of end-to-end speech recognition on mobile devices for named entities. In *ASRU*. IEEE, 23–30.

[34] Khe Chai Sim, Petr Zadrazil, and Françoise Beaufays. 2019. An Investigation Into On-device Personalization of End-to-end Automatic Speech Recognition Models. In *Interspeech*.

[35] Christophe Van Gysel, Mirko Hannemann, Ernie Pusateri, Youssef Oualil, and Ilya Oparin. 2022. Space-Efficient Representation of Entity-centric Query Language Models. In *Interspeech*.

[36] Christophe Van Gysel, Manos Tsagkias, Ernest Pusateri, and Ilya Oparin. 2020. Predicting entity popularity to improve spoken entity recognition by virtual assistants. In *SIGIR*. 1613–1616.

[37] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in Neural Information Processing Systems* 30 (2017).

[38] Dong Wang, Xiaodong Wang, and Shaohe Lv. 2019. An overview of end-to-end automatic speech recognition. *Symmetry* 11, 8 (2019), 1018.

[39] Xiaoqiang Xiao, Hong Chen, Mark Zylak, Daniela Sosa, Suma Desu, Mahesh Krishnamoorthy, Daben Liu, Matthias Paulik, and Yuchen Zhang. 2018. Geographic language models for automatic speech recognition. In *ICASSP*. IEEE, 6124–6128.

[40] Liyan Xu, Yile Gu, Jari Kolehmainen, Haidar Khan, Ankur Gandhe, Ariya Rastrow, Andreas Stolcke, and Ivan Bulyko. 2022. RescoreBERT: Discriminative

Speech Recognition Rescoring With Bert. In *ICASSP*. IEEE, 6117–6121.

[41] Mingbin Xu, Congzheng Song, Ye Tian, Neha Agrawal, Filip Granqvist, Rogier van Dalen, Xiao Zhang, Arturo Argueta, Shiyi Han, Yaqiao Deng, et al. 2022. Training Large-Vocabulary Neural Language Models by Private Federated Learning for Resource-Constrained Devices. *arXiv preprint arXiv:2207.08988* (2022).

[42] Xiaoxin Yin and Sarthak Shah. 2010. Building taxonomy of web search intents for name entity queries. In *WWW*. ACM, 1001–1010.

[43] Shiliang Zhang, Hui Jiang, Mingbin Xu, Junfeng Hou, and Li-Rong Dai. 2015. The fixed-size ordinally-forgetting encoding method for neural network language models. In *ACL*. 495–500.