

# Predicting Entity Popularity to Improve Spoken Entity Recognition by Virtual Assistants

Christophe Van Gysel  
cvangysel@apple.com  
Apple

Ernest Pusateri  
epusateri@apple.com  
Apple

Manos Tsagkias  
etsagkias@apple.com  
Apple

Ilya Oparin  
ioparin@apple.com  
Apple

## ABSTRACT

We focus on improving the effectiveness of a Virtual Assistant (VA) in recognizing emerging entities in spoken queries. We introduce a method that uses historical user interactions to forecast which entities will gain in popularity and become trending, and it subsequently integrates the predictions within the Automated Speech Recognition (ASR) component of the VA. Experiments show that our proposed approach results in a 20% relative reduction in errors on emerging entity name utterances without degrading the overall recognition quality of the system.

## KEYWORDS

virtual assistants; spoken IR; entity popularity; ASR

### ACM Reference Format:

Christophe Van Gysel, Manos Tsagkias, Ernest Pusateri, and Ilya Oparin. 2020. Predicting Entity Popularity to Improve Spoken Entity Recognition by Virtual Assistants. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '20)*, July 25–30, 2020, Virtual Event, China. ACM, New York, NY, USA, 4 pages. <https://doi.org/10.1145/3397271.3401298>

## 1 INTRODUCTION

Virtual assistants are becoming increasingly popular [14] to help users accomplish a variety of common tasks [6, 11]. Maarek [11] points out that the Automated Speech Recognition (ASR) engine – virtual assistants’ primary component to analyze spoken commands – is key to customer experience, as ASR errors have a detrimental effect on the ability of the virtual assistant to respond accurately to user queries. A prevalent and particularly challenging type of query is entity name queries [18] – queries that consist of solely an entity’s name (e.g., “Jennifer Aniston” without any other search terms). An analysis of a month of anonymized query logs from a popular virtual assistant shows that, for popular entities occurring in knowledge queries (i.e., the 20% most frequent entities that make up 91% of knowledge query traffic), the entity’s name by itself is the most frequent query form for 45% of these popular entities. Hence, entity name queries are a frequent use-case of voice search [7].

ASR systems are reasonably effective for popular entities as there is ample training data [16], however, effectiveness degrades for less popular entities. From a customer experience viewpoint, the problem is exacerbated for emerging entities that appear in the spotlight of attention, e.g., a news entity. In fact, if we consider entities whose query frequency (temporarily) increased by a factor of 3 from June to July 2019, the recognition accuracy of these emerging entities is 20% worse when compared to popular entities of the same months.

In this paper we focus on improving speech recognition accuracy for name entity queries, i.e., queries that consist of the name of an entity without any context – which typically helps recognition; e.g., “what is borrelia” gets recognized correctly, while saying “borrelia” gets misrecognized as *gorilla*. In particular, we are interested in improving queries revolving around long-tail and emerging entities. We introduce a hybrid system that brings Information Retrieval (IR) methods into ASR models. Within the ASR part, we introduce a model-agnostic approach that allows tuning the recognition likelihood of a set of given texts at runtime. We then adapt IR techniques from named entity recognition and entity linking to predict what emerging entities will become trending, i.e., entities that will be increasingly queried in our query logs. We use the outcome of our predictions to influence the output of the ASR system at runtime.

Our research questions are as follows: **(RQ1)** Can we predict entity name popularity in virtual assistant query logs and subsequently use those predictions to improve speech recognition in entity name queries? **(RQ2)** Does adding more historical data improve the accuracy of entity name popularity prediction? **(RQ3)** Do the various signals (i.e., features) perform well when used in isolation? We contribute: (1) a method that improves automatic speech recognition for emerging entity names; (2) a novel connection between ASR and IR that demonstrates how IR research can be applied to ASR problems, and (3) an analysis of the signals within the virtual assistant query log that indicates whether entities are becoming increasingly popular.

## 2 ASR PRIMER

Automated speech recognition is the task of translating a speech signal  $X$  into a string of words  $s$ . It works by optimizing [9, p. 289]:

$$\begin{aligned} s^* &= \operatorname{argmax}_s P(s | X) = \operatorname{argmax}_s \frac{P(X | s)}{P(X)} P(s) \\ &= \operatorname{argmax}_s P(X | s) P(s), \end{aligned} \quad (1)$$

SIGIR '20, July 25–30, 2020, Virtual Event, China

© 2020 Association for Computing Machinery.

This is the author’s version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '20)*, July 25–30, 2020, Virtual Event, China, <https://doi.org/10.1145/3397271.3401298>.

where  $P(X | s)$  is provided by the acoustic model and denotes the likelihood of speech signal  $X$  given the string of words  $s$ , and  $P(s)$  is provided by the language model (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 acoustic model provides probability estimates for a sequence of short audio frames (e.g., 10ms) given an input sequence of sound units (phones), and the lexicon that defines possible pronunciations for different words. The LM [1, 8] provides probability estimates for words following each other, thus taking lexical context into account. The acoustic model and lexicon (§5.1) are out of scope for this paper as we are specifically targeting entity names that were already recognized previously but fail to be recognized as stand-alone queries, and hence, fall into the domain of the LM.

The LM (§5.1) builds on the chain rule of probability:

$$P(W) = P(w_1 w_2 \dots w_n) = \prod_{i=1}^N P(w_i | w_1 w_2 \dots w_{i-1}). \quad (2)$$

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

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

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

### 3 BOOSTING ASR OF ENTITIES

Similar to search [15], speech systems must return a hypothesis quickly as users expect fast respond times. Hence, Eq. 1 is solved approximately using a greedy beam search. In addition, various components of the ASR system can be inexact. Consequently, to improve the recognition of entity names, we increase the probability of the LM to emit an entity such that the likelihood that the correct hypothesis is chosen increases. Regardless of the exact implementation of the LM conditional probability,  $P(\text{next} | \text{history})$ , we can artificially enrich the LM training data to include <s> ENTITY </s> with a predefined probability of  $\alpha$ , where ENTITY corresponds to a special word which can be substituted dynamically at runtime. More specifically, at runtime, whenever ENTITY occurs within the LM provided by  $P$ , we insert the conditional probability distribution  $Q(\cdot | \text{ENTITY})$ , which equals the predicted popularity of an entity name between the time the ASR system was built and the time it was used to handle user requests. In this paper, we select a list of potentially-popular entity names (; out of a large candidate set) and construct  $Q(\cdot | \text{ENTITY})$  by assigning an equal probability to every entity name in the list.

### 4 FORECASTING TRENDING ENTITY NAMES

We present a framework for obtaining a set of entity names that are likely to be queried more frequently. Denote  $T_i = [t_i, t_{i+1})$  as the time period that starts at time  $t_i$  and ends at time  $t_{i+1}$  (excl.). We gather entity query (§5.1) statistics over a series of  $n$  consecutive, equi-length time periods,  $T_1, \dots, T_n$ , each of length  $|T|$  (1 week in this paper). For every time period  $T_i$ , we have a set of entities  $e_j \in E^{T_i}$  with associated query frequency  $f(T_i, e_j)$ , i.e., the number of times entity  $e_j$  occurred in a query during time period  $T_i$ . The

**Table 1: Overview of time series features used in our machine-learned models. To avoid numerical issues due to extreme values, we define  $\log(0) = 10^{-6}$  and clip the individual features to a maximum value of  $10^3$ .**

	Description	Mathematical expression
$F_1$	rel. freq.	$P(e_j   T_i) = f(T_i, e_j) / \sum_k f(T_i, e_k)$
$F_2$	delta rel. freq.	$P(e_j   T_i) - P(e_j   T_{i-1})$
$F_3$	delta log rel. freq.	$\log P(e_j   T_i) - \log P(e_j   T_{i-1})$
$F_4$	rel. delta freq.	$(f(T_i, e_j) - f(T_{i-1}, e_j)) / f(T_{i-1}, e_j)$
$F_5$	rel. delta log freq.	$(\log f(T_i, e_j) - \log f(T_{i-1}, e_j)) / \log f(T_{i-1}, e_j)$
$F_6$	rel. freq w/ compl.	$P(e_j   T_i) \cdot (1.0 - P(e_j   T_{i-1}))$
$F_7$	ratio freq.	$f(T_i, e_j) / f(T_{i-1}, e_j)$

goal is now to predict the set of entities  $\mathbb{T}(E^{T_n}) \subset E_{T_1}^{T_{n-1}}$ , where  $E_{T_1}^{T_{n-1}} = \bigcup_{k=1}^{n-1} E^{T_k}$ , whose frequency will increase by a factor of  $c$  from time period  $T_{n-1}$  to time period  $T_n$ , i.e.,

$$\mathbb{T}(E^{T_n}) = \{e_j \in E_{T_1}^{T_{n-1}} \mid f(T_n, e_j) \geq c \cdot f(T_{n-1}, e_j)\}. \quad (3)$$

We cast our task as a supervised classification problem where features are extracted from  $T_1, \dots, T_{n-1}$  and the binary target labels are determined by membership in  $\mathbb{T}(E^{T_n})$ . Table 1 shows an overview of our features. The number of features varies with the number of weeks we use for feature generation (i.e.,  $n - 1$ ). Specifically, if we have  $n - 1$  weeks for feature generation, then the number of features equals  $(n - 1) + (n - 2) \cdot 6$  (i.e., one feature per week and 6 features for every pair of consecutive weeks).

## 5 EXPERIMENTAL SETUP

### 5.1 Virtual assistant system

We experiment with an ASR system consisting of statically interpolated 4-gram LMs [2, 12], the interpolated LM is entropy pruned [17]. The LMs are trained on a wide variety of data sources, such as manually and automatically transcribed virtual assistant requests and, artificially-generated sources generated from knowledge base statistics. The acoustic model is a convolutional 28M neural network trained from millions of manually transcribed, anonymized virtual assistant requests. The input to the acoustic model consists of 40 mel-spaced filter bank outputs [9]; each frame is concatenated with the preceding and succeeding ten frames, giving an input of 840 dimensions. The lexicon is manually curated by a team of experts and updated regularly with pronunciations of entities.

We set  $\alpha$  (§3) empirically by analyzing anonymized query logs from January to May 2019; and setting  $\alpha$  proportional to the number of times the full query matched a common person name. When we construct the conditional probability distribution  $Q(\cdot | \text{ENTITY})$  for which entities to boost (§3), we filter out entities that can already be correctly recognized by performing a simulation (i.e., a feedback loop) that consists of (a) automatically synthesizing voice queries for the entire list of entities using a state-of-the-art Text-To-Speech (TTS) system, (b) performing recognition using the ASR system without boosting entities, and (c) removing those entities from our list of entities for which the recognition hypothesis is an exact match of the entity name. After ASR, voice assistant queries are classified into a set of intents by an LSTM model that uses word

embedding features. We focus only on knowledge queries of encyclopedic nature that consist of an entity name with optional context [13]. Knowledge queries are parsed by matching entity names in a manually-curated knowledge base that is a superset of Wikipedia.

## 5.2 Methods under comparison

**5.2.1 Machine learning methods.** We experiment with two classifiers (§4): (a) AdaBoost [3] with 50 decision trees of depth 1, and (b) a feed-forward neural network with 2 hidden layers, 128 neurons (Glorot initialization [4]) and ReLU [5] as activation function, trained by minimizing the cross-entropy loss of the logistic activation at the output layer using Adam ( $\alpha = 10^{-3}$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 10^{-8}$ ) [10] for as many iterations until the loss converges.

**5.2.2 Baselines.** We consider five baseline systems: an ASR system without boosting entities, and four naive scoring techniques that we use to obtain a ranking of hypotheses: (a) RANDOM, a random score for every entity, (b) POPULAR LAST WEEK, the previous week’s frequency (Table 1;  $F_1$ ), (c) SUDDENLY POPULAR (Table 1;  $F_6$ ), and (d) TRENDING LAST WEEK, the ratio of the frequency of the previous week with the week before that (Table 1;  $F_7$ ).

## 5.3 Datasets and evaluation measures

We use a sub-sample of knowledge-oriented anonymized query logs of a popular voice assistant service that originated from the USA in English between June 10 and August 4, 2019. Sampling from the logs imposes an implicit absolute threshold on the frequency counts, as counts lower than the inverse sample rate are likely to be filtered out.<sup>1</sup> The sample contains over 300k unique entity names (§5.1). The threshold factor,  $c$  (Eq. 3), is set to 3. The last week of this period (July 29 to August 4) is used for evaluation (i.e., test set) where the set of trending entities is determined using Eq. 3 (5 583 out of 126 399 entities were trending). The preceding week, July 22–28, is used as target for training models (5 198 out of 125 934 entities were trending). For feature generation (i.e., the input to our models; §4), unless otherwise stated, we use three weeks of statistics (i.e., features are generated from July 1–21 for training and July 8–28 for testing) within our experiments.

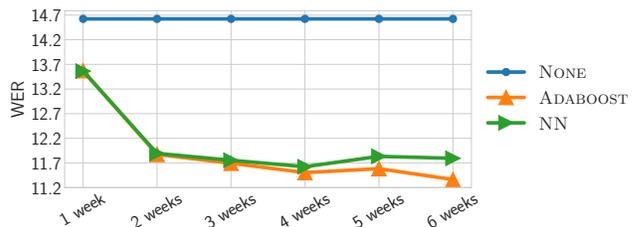
We measure the effectiveness of our approach by evaluating both the classification and ASR accuracy on the test set. More specifically, we measure: (a) the classification effectiveness, by computing average precision (AP), recall@ $k$ , precision@ $k$  for  $k = 2\,500, 5\,000, 10\,000$  where  $k$  is the position at which we cut-off the ranking, and (b) the ASR accuracy, by computing word error rate (WER) where we boost the top- $k$  (i.e., WER@ $k$ ; selected after applying the feedback loop described in §5.1) entity names using the approach described in Section §3 while performing ASR on the automatically-synthesized (using a state-of-the-art TTS system; §5.1) names of trending entities in our test set (§5.3). Hyperparameter  $k$  is our budget of entity names we are able to boost at runtime; a value of  $k$  beyond 10 000 could lead to increased latency in handling user requests. Statistical significance of observed differences in per-utterance WER is determined using a two-tailed paired Student’s t-test (\*  $p < 0.01$ ).

## 6 RESULTS AND DISCUSSION

In this section, we answer the three research questions posed in §1.

<sup>1</sup>We are unable to provide the exact sample rate due to confidentiality; however, the sample rate is less than 15%.

**RQ1:** Table 2 shows a comparison of the ASR system without any entity name boosting (WITHOUT BOOSTING), four heuristic methods (§5.2.2) that generate a ranking of entity names to boost, and the machine learning models (§5.2.1) trained on the features outlined in Table 1 with various rank cut-offs  $k$  (§5.3). None of the heuristics (§5.2.2) contribute a meaningful improvement. In fact, most of the heuristics actually degrade WER ( $k = 2\,500$ ) on our test set of spoken entity names (§5.3). For the machine-learned models (ADABOOST, NN; §5.2.1), we see that both models improve WER significantly between 10% and 20% (depending on the value of  $k$ ) compared to the system without any entity name boosting. A qualitative investigation of the predictions shows that the heuristics perform poorly as they rank popular entities highly; as opposed to infrequent entities that are emerging. For example, all heuristics (except RANDOM) rank the entity name *Robert Mueller* highest. On the other hand, the rankings predicted by the machine-learned models have *Scripture*, *Six Sigma*, and *duchess* amongst their top-ranked entity names. We will revisit this observation when we discuss **RQ3**. Finally, we used the systems in the bottom two rows of Table 2 (i.e., ADABOOST, NN) to recognize a test set consisting of a uniform random sample of manually-annotated, general virtual assistant queries (as opposed to entity-bearing knowledge queries) and measured only a negligible degradation of WER (less than 0.01% relative; not shown in table). Hence, we believe our approach can be used without negatively impacting the user experience.



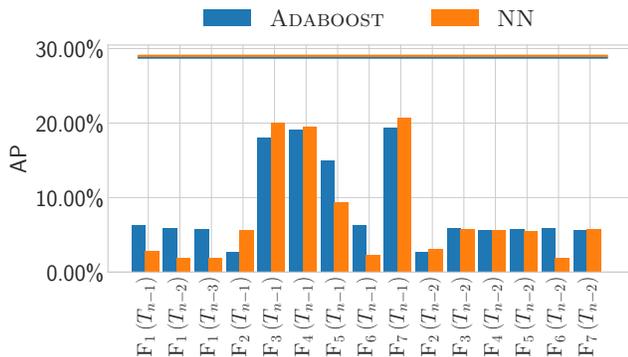
**Figure 1: The effect of varying the amount of time periods (in weeks) to extract features from; with  $k = 5\,000$ .**

**RQ2:** Figure 1 shows the behavior of the WER test metric as we include an increasing amount of history in our models. We can see that even with a single week of history, i.e., models with only feature  $F_1(n-1)$  as input, we obtain a 7% reduction in WER. Including additional weeks, and hence, rate-of-change features, further reduces WER to around the same level of effectiveness as reported in Table 2 (i.e., 3 weeks of history) and eventually WER reaches a plateau around 5 to 6 weeks of history. Hence, as expected, recent history is more important than far-away history for the entity popularity forecasting task. However, Figure 1 raises the question regarding why a model with only feature  $F_1(n-1)$  as input performs better than the heuristic POPULAR LAST WEEK? We answer this question together with **RQ3**.

**RQ3:** Figure 2 shows the classification effectiveness (AP) when we train a model for every feature individually ( $n = 4$ ) using our machine learning models (§5.2.1). Figure 2 confirms our findings in **RQ2**, namely that rate-of-change features play an important role. However, we note that none of the features individually come close to the same level of AP as was obtained in Table 2 by combining all features. In addition, many of our features are correlated

**Table 2: Overview of experimental results (all metrics are in percentages). We generate a ranked list of entities according to predicted popularity and boost the top- $k$  ( $k = 2500, \dots$ ) entity names at recognition time. Machine-learned models that use historical entity query interaction statistics can significantly (\*; §5.3) reduce the WER (best in bold) for trending entities.**

	2500				5000			10000		
	AP	Recall	Precision	WER	Recall	Precision	WER	Recall	Precision	WER
WITHOUT BOOSTING	-	-	-	14.67	-	-	14.67	-	-	14.67
<b>Heuristics</b>										
RANDOM	2.57	1.00	2.24	14.50	2.06	2.30	14.50	4.35	2.43	14.18
POPULAR LAST WEEK	3.73	0.16	0.36	14.86	0.38	0.42	14.95	0.97	0.54	14.58
SUDDENLY POPULAR	2.80	0.16	0.36	14.85	0.38	0.42	14.97	0.97	0.54	14.64
TRENDING LAST WEEK	3.78	0.21	0.48	14.88	0.66	0.74	14.69	1.65	0.92	13.56
<b>Machine-learned</b>										
ADABOOST	28.79	20.31	45.36	<b>12.29*</b>	33.12	36.98	<b>11.74*</b>	49.02	27.37	11.75*
NN	29.02	20.85	46.56	12.30*	32.40	36.18	11.80*	48.56	27.11	<b>11.72*</b>



**Figure 2: Effectiveness of models trained on individual features in terms of AP; with  $k = 5000$ .**

(e.g.,  $F_3$  is the same as  $F_2$  with the frequencies in log-space), and hence, the effectiveness of individual features is not additive. Consequently, we conclude from Figure 2 that combining our different features is useful for the entity popularity forecasting task. Lastly, we revisit the observation raised in **RQ1**, where we noted that the heuristics (§5.2.2) are biased towards head entities. Most notably, heuristic TRENDING LAST WEEK obtains an AP of about 4% in Table 2, whereas the AP of the models trained with TRENDING LAST WEEK (i.e.,  $F_7(T_{n-1})$ ) in Figure 2 is close to 20%. The difference here is due to the fact that our machine learning models (§5.2.1) learn to threshold feature values (i.e., the decision trees in ADABOOST and ReLU activation in NN). Hence, the feature transformations learned by our models are used to filter out head entities and this model property is key to their effectiveness.

## 7 CONCLUSIONS

We introduced a method to improve the speech recognition quality of queries revolving around emerging entities. We trained a model to forecast what entities will become popular from historical entity popularity data, and we subsequently used that information to boost the recognition likelihood of the predicted trending entities. A comparison of our method with five baselines showed solid improvements of 20% in word error rate without degrading performance on a test set of general assistant voice queries. Our analyses showed that (a) increasing the amount of history—beyond 1 week—used to compute features helps to some extent, (b) our features are complementary, and (c) non-linear feature transformations are important for our features to be effective. Experimental

results not included in this paper show that our findings are consistent over multiple different time periods. The primary limitation of our work is that we only make use of information contained within the virtual assistant query log. Hence, our approach is limited only to entities that can be recognized already, possibly preceded by additional spoken context that disambiguates the utterance. Future work can focus on out-of-vocabulary entities and entities that are often misrecognized, even when additional spoken context is available. In addition, external data sources could provide a useful signal that is not available within the virtual assistant query logs.

**Acknowledgments.** The authors would like to thank Russ Webb, Barry Theobald, Rolf Jagerman, the anonymous, and the internal reviewers for their valuable comments.

## REFERENCES

- [1] Y. Bengio, R. Ducharme, P. Vincent, and C. Janvin. A neural probabilistic language model. *JMLR*, 3:1137–1155, 2003.
- [2] K. W. Church and W. A. Gale. A comparison of the enhanced Good-Turing and deleted estimation methods for estimating probabilities of English bigrams. *Computer Speech & Language*, 5(1):19–54, 1991.
- [3] Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, 55(1):119–139, 1997.
- [4] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, pages 249–256, 2010.
- [5] X. Glorot, A. Bordes, and Y. Bengio. Deep sparse rectifier neural networks. In *AISTATS*, pages 315–323, 2011.
- [6] D. Graus, P. N. Bennett, R. W. White, and E. Horvitz. Analyzing and predicting task reminders. In *UMAP*, pages 7–15, 2016.
- [7] I. Guy. The characteristics of voice search: comparing spoken with typed-in mobile web search queries. *TOIS*, 36(3):30, 2018.
- [8] F. Jelinek. *Statistical Methods for Speech Recognition*. MIT Press, Cambridge, MA, USA, 1998.
- [9] D. Jurafsky and J. H. Martin. *Speech and Language Processing, 2nd edition*. Prentice Hall, 2008.
- [10] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization, 2014.
- [11] Y. Maarek. Alexa, can you help me shop? In *SIGIR*, pages 1369–1370, 2019.
- [12] E. Pusateri, C. Van Gysel, R. Botros, S. Badaskar, M. Hannemann, Y. Oualil, and I. Oparin. Connecting and comparing language model interpolation techniques. In *Interspeech*, pages 3500–3504, 2019.
- [13] R. Reinanda, E. Meij, and M. de Rijke. Mining, ranking and recommending entity aspects. In *SIGIR*, pages 263–272. ACM, 2015.
- [14] J. Research. Digital voice assistants in use to triple to 8 billion by 2023, driven by smart home devices. Press Release, Feb. 2019.
- [15] E. Schurman and J. Brutlag. Performance related changes and their user impact. Presentation at Velocity – Web Performance and Operations Conf., 2009.
- [16] J. Serrino, L. Velikovich, P. Aleksic, and C. Allauzen. Contextual Recovery of Out-of-Lattice Named Entities in Automatic Speech Recognition. In *Interspeech*, pages 3830–3834, 2019.
- [17] A. Stolcke. Entropy-based pruning of backoff language models. In *DARPA Broadcast News Transcription and Understanding Workshop*, pages 270–274, 2000.
- [18] X. Yin and S. Shah. Building taxonomy of web search intents for name entity queries. In *WWW*, pages 1001–1010. ACM, 2010.