# QnAMaker: Data to Bot in 2 Minutes

Parag Agrawal
Microsoft Corporation
paragag@microsoft.com

Tulasi Menon
Microsoft Corporation
tulasim@microsoft.com

Aya Kamel
Microsoft Corporation
aykame@microsoft.com

Michel Naim
Microsoft Corporation
mgerguis@microsoft.com

Chaikesh Chouragade
Microsoft Corporation
chchoura@microsoft.com

Gurvinder Singh
Microsoft Corporation
gurvsing@microsoft.com

Rohan Kulkarni
Microsoft Corporation
rokulka@microsoft.com

Anshuman Suri
Microsoft Corporation
ansuri@microsoft.com

Sahithi Katakam
Microsoft Corporation
sakataka@microsoft.com

Vineet Pratik
Microsoft Corporation
vipratik@microsoft.com

Prakul Bansal
Microsoft Corporation
prabansa@microsoft.com

Simerpreet Kaur
Microsoft Corporation
sikaur@microsoft.com

Neha Rajput
Microsoft Corporation
nerajput@microsoft.com

Anand Duggal
Microsoft Corporation
anduggal@microsoft.com

Achraf Chalabi
Microsoft Corporation
achalabi@microsoft.com

Prashant Choudhari
Microsoft Corporation
pchoudh@microsoft.com

Somi Reddy Satti
Microsoft Corporation
sosatti@microsoft.com

Niranjan Nayak
Microsoft Corporation
niranjan@microsoft.com

## ABSTRACT

Having a bot for seamless conversations is a much-desired feature that products and services today seek for their websites and mobile apps. These bots help reduce traffic received by human support significantly by handling frequent and directly answerable known questions. Many such services have huge reference documents such as FAQ pages, which makes it hard for users to browse through this data. A conversation layer over such raw data can lower traffic to human support by a great margin. We demonstrate QnAMaker, a service that creates a conversational layer over semi-structured data such as FAQ pages, product manuals, and support documents. QnA-Maker is the popular choice for Extraction and Question-Answering as a service and is used by over 15,000 bots in production. It is also used by search interfaces and not just bots.

## KEYWORDS

ChatBots, Democratizing AI, Question Answering, Online learning, Bot Persona, Information Retrieval, Multilingual

## 1 INTRODUCTION

QnAMaker aims to simplify the process of bot creation by extracting Question-Answer (QA) pairs from data given by users into a Knowledge Base (KB) and providing a conversational layer over it. KB here refers to one instance of azure search index, where the extracted QA are stored. Whenever a developer creates a KB using QnAMaker, they automatically get all NLP capabilities required to answer user's queries. There are other systems such as Google's Dialogflow, IBM's Watson Discovery which tries to solve this problem. QnAMaker provides unique features for the ease of development such as the ability to add a persona-based chit-chat layer on top of the bot. Additionally, bot developers get automatic feedback from the system based on end-user traffic and interaction which helps them in enriching the KB; we call this feature active-learning[1]. Our system also allows user to add Multi-Turn structure to KB using hierarchical extraction and contextual ranking. QnAMaker today supports over 35 languages, and is the only system among its competitors to follow a Server-Client architecture; all the KB data rests only in the client's subscription, giving users total control over

---

[1] https://docs.microsoft.com/en-us/azure/cognitive-services/qnamaker/how-to/improve-knowledge-base
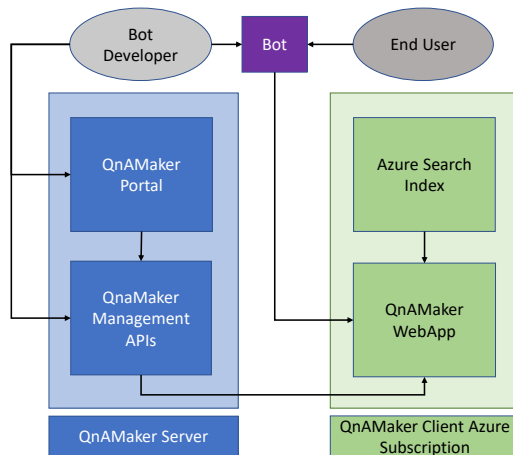
**Figure 1: Interactions between various components of Qna-Maker, along with their scopes: server-side and client-side**

their data. QnAMaker is part of Microsoft Cognitive Service and currently runs using the Microsoft Azure Stack[2].

## 2 SYSTEM DESCRIPTION

### 2.1 Architecture

As shown in Figure 1, humans can have two different kinds of roles in the system: **Bot-Developers** who want to create a bot using the data they have, and **End-Users** who will chat with the bot(s) created by bot-developers. The components involved in the process are:

- *QnAMaker Portal*[3]: This is the Graphical User Interface (GUI) for using QnAMaker. This website is designed to ease the use of management APIs. It also provides a test pane.
- *QnaMaker Management APIs*: This is used for the extraction of Question-Answer (QA) pairs from semi-structured content. It then passes these QA pairs to the web app to create the Knowledge Base Index.
- *Azure Search Index*: Stores the KB with questions and answers as indexable columns, thus acting as a retrieval layer.
- *QnaMaker WebApp*: Acts as a layer between the Bot, Management APIs, and Azure Search Index. WebApp does ranking on top of retrieved results. WebApp also handles feedback management for active learning.
- *Bot*: Calls the WebApp with the User's query to get results.

### 2.2 Bot Development Process

Creating a bot is a 3-step process for a bot developer:

(1) Create a QnaMaker Resource in Azure: This creates a WebApp with binaries required to run QnAMaker. It also creates an Azure Search Service for populating the index with any given knowledge base, extracted from user data

(2) Use Management APIs to Create/Update/Delete your KB: The Create API automatically extracts the QA pairs and

sends the Content to WebApp, which indexes it in Azure Search Index. Developers can also add persona-based chat content and synonyms while creating and updating their KBs.

(3) Bot Creation: Create a bot using any framework and call the WebApp hosted in Azure to get your queries answered. There are Bot-Framework templates [4] provided for the same.

### 2.3 Extraction

The Extraction component is responsible for understanding a given document and extracting potential QA pairs. These QA pairs are in turn used to create a KB to be consumed later on by the *QnA-Maker WebApp* to answer user queries. First, the basic blocks from given documents such as text, lines are extracted. Then the layout of the document such as columns, tables, lists, paragraphs, etc is extracted. This is done using Recursive X-Y cut [4]. Following Layout Understanding, each element is tagged as headers, footers, table of content, index, watermark, table, image, table caption, image caption, heading, heading level, and answers. Agglomerative clustering [3] is used to identify heading and hierarchy to form an intent tree. Leaf nodes from the hierarchy are considered as QA pairs. In the end, the intent tree is further augmented with entities using CRF-based sequence labeling. Intents that are repeated in and across documents are further augmented with their parent intent, adding more context to resolve potential ambiguity.

### 2.4 Retrieval And Ranking

QnAMaker uses *Azure Search Index* as it's retrieval layer, followed by re-ranking on top of retrieved results (Figure 2). Azure Search is based on inverted indexing and TF-IDF scores. Azure Search provides fuzzy matching based on edit-distance, thus making retrieval robust to spelling mistakes. It also incorporates lemmatization and normalization. These indexes can scale up to millions of documents, lowering the burden on *QnAMaker WebApp* which gets less than 100 results to re-rank.

Different customers may use QnAMaker for different scenarios such as banking task completion, answering FAQs on company policies, or fun and engagement. The number of QAs, length of questions and answers, number of alternate questions per QA can vary significantly across different types of content. Thus, the ranker model needs to use features that are generic enough to be relevant across all use cases.

*2.4.1 Pre-Processing.* The pre-processing layer uses components such as Language Detection, Lemmatization, Speller, and Word Breaker to normalize user queries. It also removes junk characters and stop-words from the user's query.

*2.4.2 Features.* Going into granular features and the exact empirical formulas used is out of the scope of this paper. The broad level features used while ranking are:

(1) **WordNet:** There are various features generated using WordNet [8] matching with questions and answers. This takes care of word-level semantics. For instance, if there is information about "price of furniture" in a KB and the end-user asks about "price of table", the user will likely get a relevant

| Domain | Number of QAs | Avg Questions per QA | AUC (%) | $F_1$ (top answer) |
|---|---|---|---|---|
| Navigation Help Bot | 56 | 12.5 | 88.7 | 71.2 |
| Chit-Chat Alone | 100 | 9.8 | 92.4 | 88.6 |
| CustomerCare Interface | 164 | 2.2 | 90.9 | 86.7 |
| HR Internal Bot | 52 | 1.0 | 85.5 | 82.6 |
| HR Internal Bot (with Chit-Chat) | 152 | 6.78 | 82.7 | 77.6 |

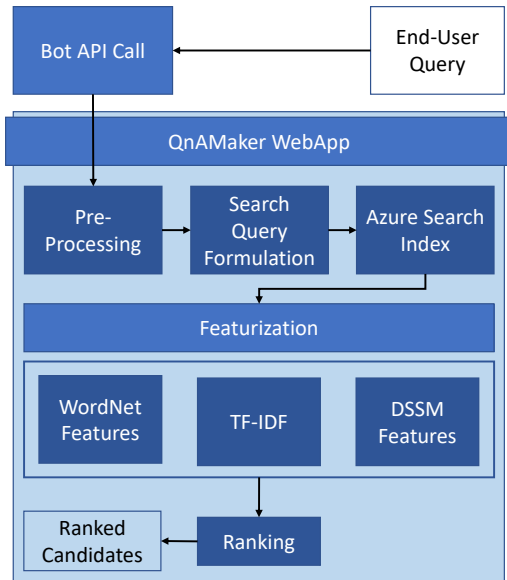**Table 1: Retrieval And Ranking Measurements**



**Figure 2: QnAMaker Runtime Pipeline**

answer. The scores of these WordNet features are calculated as a function of:

- Distance of 2 words in the WordNet graph
- Distance of Lowest Common Hypernym from the root
- Knowledge-Base word importance (Local IDFs)
- Global word importance (Global IDFs)

This is the most important feature in our model as it has the highest relative feature gain.

(2) **CDSSM:** Convolutional Deep Structured Semantic Models [6] are used for sentence-level semantic matching. This is a dual encoder model that converts text strings (sentences, queries, predicates, entity mentions, etc) into their vector representations. These models are trained using millions of Bing Query Title Click-Through data. Using the *source-model* for vectorizing user query and *target-model* for vectorizing answer, we compute the cosine similarity between these two vectors, giving the relevance of answer corresponding to the query.

(3) **TF-IDF:** Though sentence-to-vector models are trained on huge datasets, they fail to effectively disambiguate KB specific data. This is where a standard TF-IDF [7] featurizer with local and global IDFs helps.

*2.4.3 Contextual Features.* We extend the features for contextual ranking by modifying the candidate QAs and user query in these ways:

- $Query_{modified}$ = Query + Previous Answer; For instance, if user query is "yes" and the previous answer is "do you want to know about XYZ", the current query becomes "do you want to know about XYZ yes".
- Candidate QnA pairs are appended with its parent Questions and Answers; no contextual information is used from the user's query. For instance, if a candidate QnA has a question "benefits" and its parent question was "know about XYZ", the candidate QA's question is changed to "know about XYZ benefits".

The features mentioned in Section 2.4.2 are calculated for the above combinations also. These features carry contextual information.

*2.4.4 Modeling and Training.* We use gradient-boosted decision trees as our ranking model to combine all the features. Early stopping [11] based on Generality-to-Progress ratio is used to decide the number of step trees and Tolerant Pruning [10] helps prevent overfitting. We follow incremental training if there is small changes in features or training data so that the score distribution is not changed drastically.

## 2.5 Persona Based Chit-Chat

We add support for bot-developers to directly enable handling chit-chat queries like "hi", "thank you", "what's up" in their QnAMaker bots. In addition to chit-chat, we also give bot developers the flexibility to ground responses for such queries in a specific personality: professional, witty, friendly, caring, or enthusiastic. For example, the "Humorous" personality can be used for a casual bot, whereas a "Professional" personality is more suited in case of banking FAQs or task-completion bots. There is a list of 100+ predefined intents [1]. There is a curated list of queries for each of these intents, along with a separate query understanding layer for ranking these intents. The arbitration between chit-chat answers and user's knowledge base answers is handled by using a chat-domain classifier [2].

## 2.6 Active Learning

The majority of the KBs are created using existing FAQ pages or manuals but to improve the quality it requires effort from the developers. Active learning generates suggestions based on end-user feedback as well as ranker's implicit signals. For instance, if for a query, CDSSM feature was confident that one QnA should be ranked higher whereas wordnet feature thought other QnA should

be ranked higher, active learning system will try to disambiguate it by showing this as a suggestion to the bot developer. To avoid showing similar suggestions to developers, DB-Scan clustering is done which optimizes the number of suggestions shown.

## 3 EVALUATION AND INSIGHTS

QnAMaker is not domain-specific and can be used for any type of data. To support this claim, we measure our system's performance for datasets across various domains. The evaluations are done by managed judges who understands the knowledge base and then judge user queries relevance to the QA pairs (binary labels). Each query-QA pair is judged by two judges. We filter out data for which judges do not agree on the label. Chit-chat in itself can be considered as a domain. Thus, we evaluate performance on given KB both with and without chit-chat data (last two rows in Table 1), as well as performance on just chit-chat data (2nd row in Table 1). Hybrid of deep learning(CDSSM) and machine learning features give our ranking model low computation cost, high explainability and significant F1/AUC score. Based on QnAMaker usage, we observed these trends:

- Around 27% of the knowledge bases created use pre-built persona-based chitchat, out of which, ~4% of the knowledge bases are created for chit-chat alone. The highest used personality is Professional which is used in 9% knowledge bases.
- Around ~25% developers have enabled active learning suggestions. The acceptance to reject ratio for active learning suggestions is 0.31.
- 25.5% of the knowledge bases use one URL as a source while creation. ~41% of the knowledge bases created use different sources like multiple URLs. 15.19% of the knowledge bases use both URL and editorial content as sources. Rest use just editorial content.

## 4 DEMONSTRATION

We demonstrate QnAMaker: a service to add a conversational layer over semi-structured user data. In addition to query-answering, we support novel features like personality-grounded chit-chat, active learning based on user-interaction feedback (Figure 3), and hierarchical extraction for multi-turn conversations (Figure 4). The goal of the demonstration will be to show how easy it is to create an intelligent bot using QnAMaker. All the demonstrations will be done on the production website [5] Demo Video can be seen here. [6]

## 5 FUTURE WORK

The system currently doesn't highlight the answer span and does not generate answers taking the KB as grounding. We will be soon supporting Answer Span [5] and KB-grounded response generation [9] in QnAMaker. We are also working on user-defined personas for chit-chat (automatically learned from user-documents). We aim to enhance our extraction to be able to work for any unstructured document as well as images. We are also experimenting on improving our ranking system by using semantic vector-based search as our retrieval and transformer-based models for re-ranking.
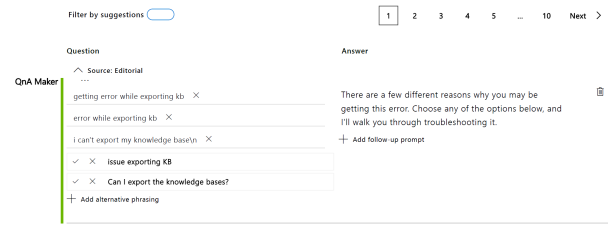
---

[5] http://www.qnamaker.ai
[6] https://youtu.be/nBmBpsjuDOo



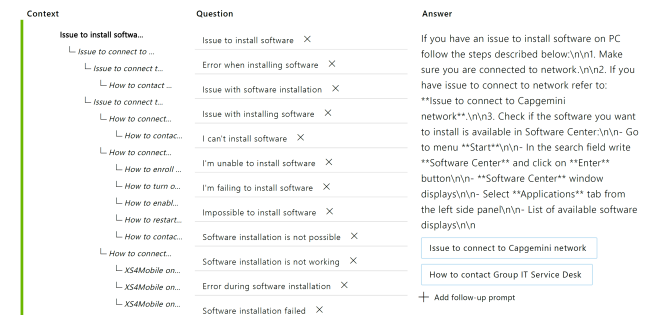**Figure 3: Active Learning Suggestions**



**Figure 4: Multi-Turn Knowledge Base**

## REFERENCES

[1] Parag Agrawal, Anshuman Suri, and Tulasi Menon. 2018. A Trustworthy, Responsible and Interpretable System to Handle Chit Chat in Conversational Bots. *CoRR* abs/1811.07600 (2018). arXiv:1811.07600 http://arxiv.org/abs/1811.07600

[2] Satoshi Akasaki and Nobuhiro Kaji. 2017. Chat Detection in an Intelligent Assistant: Combining Task-oriented and Non-task-oriented Spoken Dialogue Systems. In *ACL*.

[3] Douglass R Cutting, David R Karger, Jan O Pedersen, and John W Tukey. 1992. Scatter/gather: A cluster-based approach to browsing large document collections. In *Proceedings of the 15th annual international ACM SIGIR conference on Research and development in information retrieval.* ACM, 318–329.

[4] Jaekyu Ha, Robert M Haralick, and Ihsin T Phillips. 1995. Recursive XY cut using bounding boxes of connected components. In *Proceedings of 3rd International Conference on Document Analysis and Recognition*, Vol. 2. IEEE, 952–955.

[5] Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching Machines to Read and Comprehend. In *Advances in Neural Information Processing Systems 28*, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett (Eds.). Curran Associates, Inc., 1693–1701. http://papers.nips.cc/paper/5945-teaching-machines-to-read-and-comprehend.pdf

[6] Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. 2013. Learning deep structured semantic models for web search using clickthrough data. In *Proceedings of the 22nd ACM international conference on Conference on information & knowledge management.* ACM, 2333–2338.

[7] K Sparck Jones, Steve Walker, and Stephen E. Robertson. 2000. A probabilistic model of information retrieval: development and comparative experiments: Part 2. *Information processing & management* 36, 6 (2000), 809–840.

[8] George A. Miller. 1995. WordNet: A Lexical Database for English. *Commun. ACM* 38, 11 (Nov. 1995), 39–41. https://doi.org/10.1145/219717.219748

[9] Lianhui Qin, Michel Galley, Chris Brockett, Xiaodong Liu, Xiang Gao, Bill Dolan, Yejin Choi, and Jianfeng Gao. 2019. Conversing by Reading: Contentful Neural Conversation with On-demand Machine Reading. In *Proc. of ACL.* https://www.microsoft.com/en-us/research/publication/conversing-by-reading-contentful-neural-conversation-with-on-demand-machine-reading/

[10] Christino Tamon and Jie Xiang. 2000. On the boosting pruning problem. In *European conference on machine learning.* Springer, 404–412.

[11] Tong Zhang, Bin Yu, et al. 2005. Boosting with early stopping: Convergence and consistency. *The Annals of Statistics* 33, 4 (2005), 1538–1579.