# MeetingQA: Extractive Question-Answering on Meeting Transcripts

Archiki Prasad[1], Trung Bui[2], Seunghyun Yoon[2],
Hanieh Deilamsalehy[2], Franck Dernoncourt[2], Mohit Bansal[1]
[1]University of North Carolina at Chapel Hill, [2]Adobe Research

ACL 2023

## Motivation & Contributions

- Abundance of meeting transcripts – typically domain-specific & information-rich long documents for building NLP systems.
- Prior works focus on summarization and extracting action items, underutilizing significant QA components of meetings.
- We introduce **MeetingQA**, an extractive QA dataset comprising questions asked by participants during a meeting and corresponding answer spans from relevant discussions.
- Questions asked by participants in MeetingQA are **longer**, **open-ended**, and **discussion-seeking** including interesting scenarios such as **rhetorical questions**, **multi-span answers** and/or answers contributed by **multiple speakers**.
- Despite high human performance (F1=84.6), the best QA models yield F1 of 57.3 making MeetingQA a challenging dataset with substantial room for improvement.



**Speaker 4:** For this whole discussion who among us is doing stuff that happens online and who's doing stuff that happens offline?

**Speaker 2:** *The basic word importance is offline as well. The combined measure might not be if we want to wait for the user to type.*

**Speaker 4:** Yeah. Okay, okay.

**Speaker 3:** *Mine's gonna be mostly using the offline. But the actual stuff it's doing will be online.* But it won't be very processor intensive or memory intensive.

**Speaker 0:** *I don't know about the search functionality, that might be online.* Depends on how its gonna work.

**Speaker 4:** *That means that at least we don't have the type of situation where somebody has to do a billion calculations on data online.* Cause that would make it a lot more like that would mean that our interface for the data would have to be a lot more careful about how it performs and and everything. *And nobody is modifying that data at at online time at all it seems. Nobody's making any changes to the actual data online.*

**Speaker 3:** Don't think so.
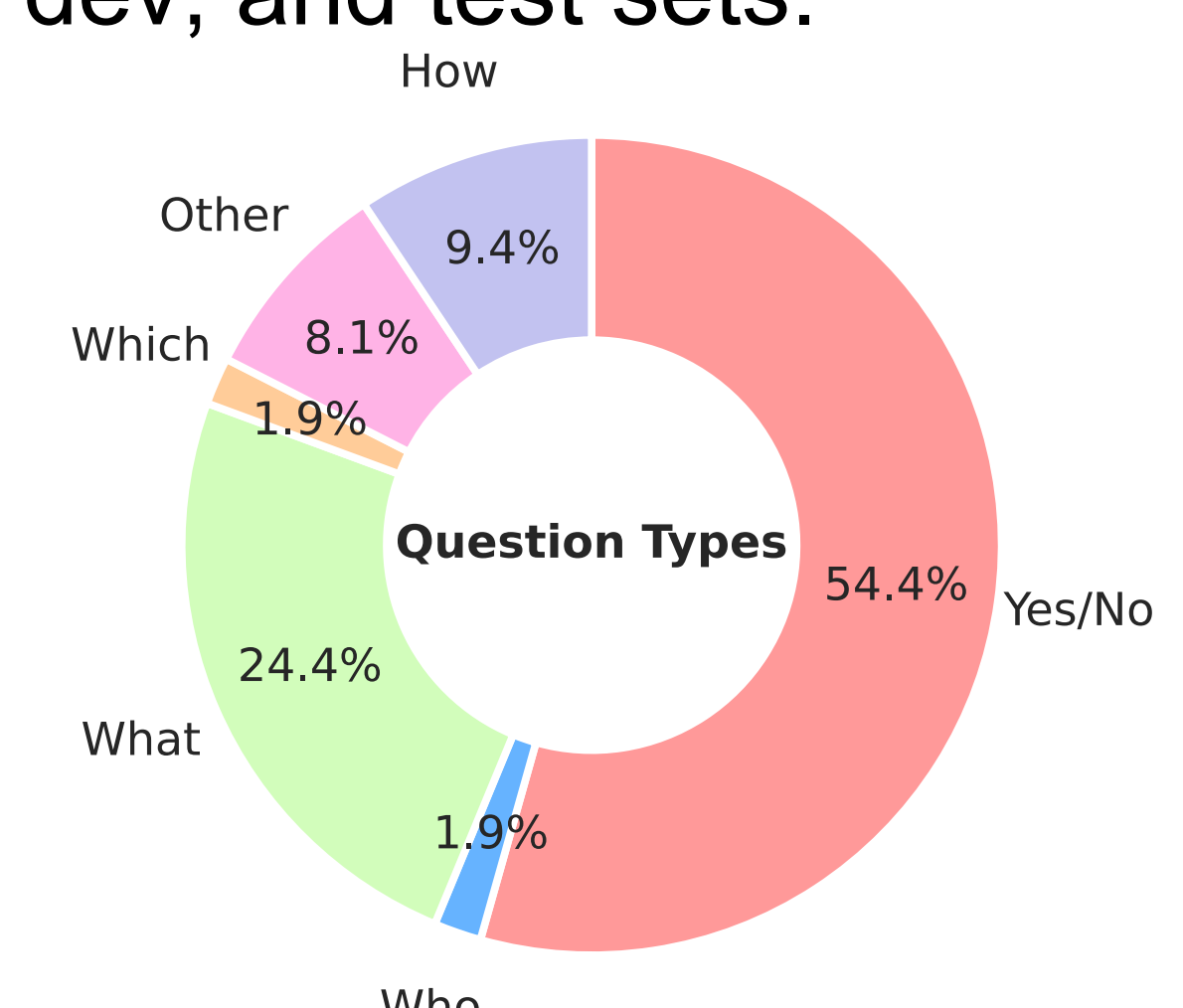
## MeetingQA: Details and Analysis

### Data Collection and Annotation

- Annotate public meetings from AMI (Augmented Multi-party Interaction) corpus with ~100 hours manually transcribed meetings.
- Question Selection: based on punctuation and number of words.
- Answer Annotation: Recruit annotators to label which sentences from the transcript answer each question along with meta-data.
- High inter-annotator agreement with Krippendorff's α of 0.73, obtaining annotations for 166 meetings at $61 per meeting.

### Dataset Information and Analysis

- Total of 7,735 questions split across train, dev, and test sets.

| | Train | Dev | Test |
|---|---|---|---|
| Number of Meetings | 64 | 48 | 54 |
| Number of Questions | 3007 | 2252 | 2476 |
| w/ No Answer | 956 | 621 | 764 |
| w/ Multi-Span Answers | 787 | 548 | 663 |
| w/ Multi-Speaker Answers | 1016 | 737 | 840 |
| Avg. Questions per Meeting | 46.98 | 46.92 | 45.85 |



Question Types: Yes/No 54.4%, What 24.4%, How 9.4%, Other 8.1%, Which 1.9%, Who 1.9%

- **Question Types**: Even questions framed in 'yes/no' manner are information-seeking and elicit detailed responses, ~50% of questions are opinion-seeking and ~20% are framed rhetorically.
- **Answer Types**: 30% of questions are unanswerable, 40% of answers are multi-span (non-consecutive sentences) and 48% involve multiple speakers. Nearly 70% of multi-speaker answers contain some level of disagreement among participants.
- **Length Distribution**: Average length of a transcript, question, and corresponding answer is 5.9K, 12, and 35 words, respectively.
- **Human Performance**: F1=84.6 on 250 questions from the test set.

## Methods and Experimental Results

### Finetuned Performance

| Model | Overall F1 | No Ans. F1 | Answerable F1 | | |
|---|---|---|---|---|---|
| | | | All | M-Span | M-Speaker |
| SS RoBERTa-base | **56.5** | 41.0 | **63.1** | **60.8** | **64.1** |
| SS Longformer-base | 55.6 | **46.1** | 59.9 | 55.3 | 59.4 |
| MS RoBERTa-base | 54.0 | 41.1 | 59.8 | 58.2 | 60.9 |
| MS Longformer-base | 53.8 | 39.4 | 60.3 | 58.8 | 62.0 |
| Human Performance | **84.6** | 80.7 | 86.3 | 88.1 | 87.7 |

- **Short-context models**: Context (that fits in model's input) retrieved from transcript based on location of the question.
- **Single-span models:** Predict single-span from first to last relevant sentence.
- **Multi-span models:** QA as token-classification.

**Results:**
(i) Short-context models slightly outperform long-context models by 1-2 F1 points.
(ii) Multi-span models have comparable or less performance than single-span models.
(iii) ≥ 25 F1 points gap with human performance.
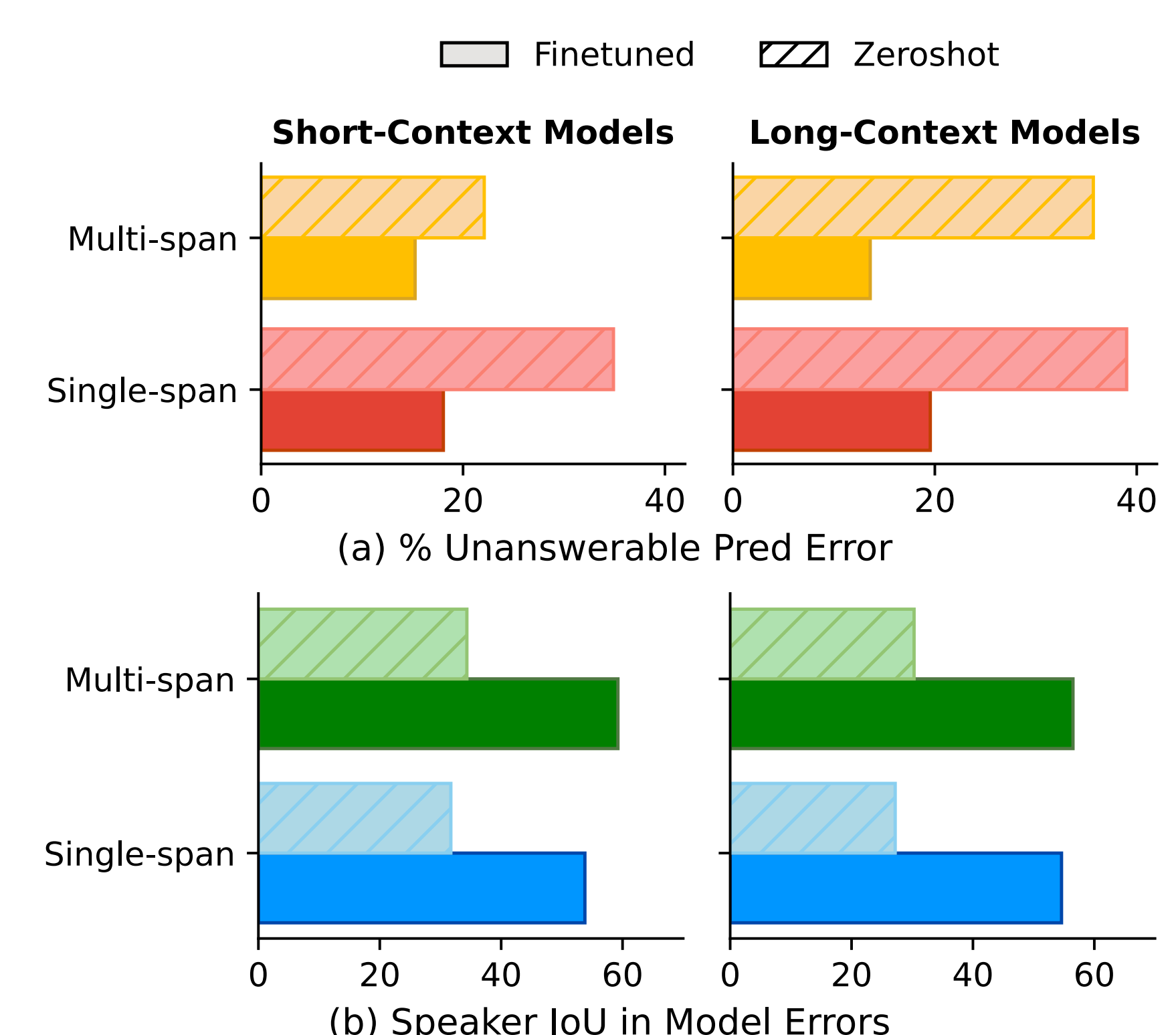
### Zero-shot Performance

| Model | Inter. Data | Overall F1 |
|---|---|---|
| RoBERTa-base | SQuADv2 | 27.9 |
| | + silver | **34.6** |
| Longformer-base | SQuADv2 | 15.1 |
| | + silver | 32.5 |
| FLAN-T5 XL | — | 33.8 |
| FLAN-T5 XL (self ans) | — | 34.0 |
| Human performance | — | **84.6** |

- **Silver Data Augmentation**: Augment training data with automatically annotated answer spans for interviews from MediaSum dataset.

**Results:**
(i) All models exhibit poor zero-shot performance (~50 F1 point gap).
(ii) Augmenting with silver data improves zero-shot performance.
(iii) Larger instruction-tuned LMs (Flan-T5) yield comparable performance.

### Error Analysis



(a) % Unanswerable Pred Error

(b) Speaker IoU in Model Errors

- Models struggle at identifying rhetorical questions, especially in zero-shot.
- Single-span predictions contain a greater fraction of irrelevant sentences.
- Models struggle to identify which speakers answer a question, especially in zero-shot setting.

*Project Page →*
*Contact: archiki@cs.unc.edu*