

# Evidence Contextualization and Counterfactual Attribution For ConvQA over Heterogeneous Data with RAG Systems

Rishiraj Saha Roy, Joel Schlotthauer, Chris Hinze,  
Andreas Foltyn, Luzian Hahn, and Fabian Küch

{rishiraj.saha.roy, joel.Schlotthauer, chris.hinze, andreas.foltyn, luzian.hahn, fabian.kuech}@iis.fraunhofer.de

## RAG over enterprise wikis is challenging

### 2024-10-02 Meeting Notes

Today we will talk about the progress of the project on retrieval augmented generation.

#### Agenda

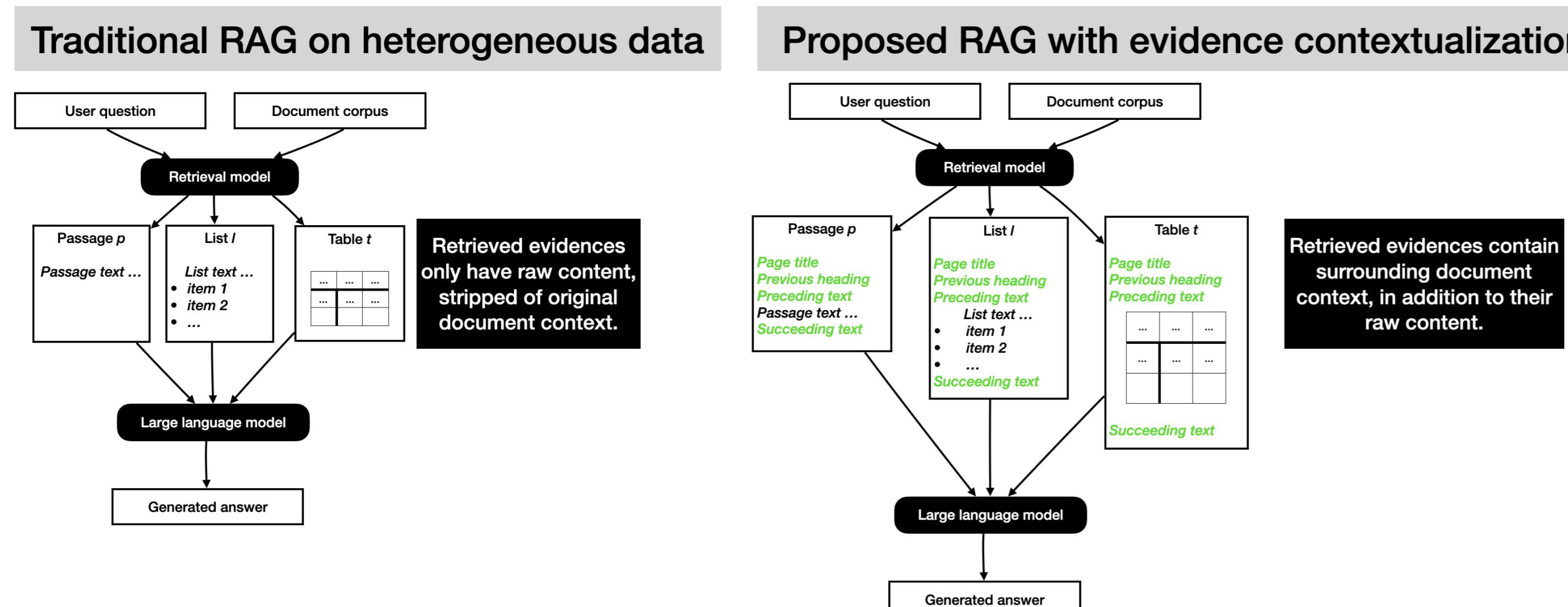
- We'll first do a basic round of RAG team updates in this month's meeting
- Everyone will report what has been done, and the to-dos

Member	Task	Action items	Time needed	Notes
Bob	Basic FE and BE	Follow-up q in UI	3 days	Currently manual
Alice	Similarity function	Fine-tune with gpt4o*	1 week	Now w/ embed cos
Trudy	Verbalizations	Batch configs*	6 hours	Running superbly

\* Alice and Trudy to fix long-standing embedding error with openxt strings

**Question:** todo for alice in oct rag meeting?

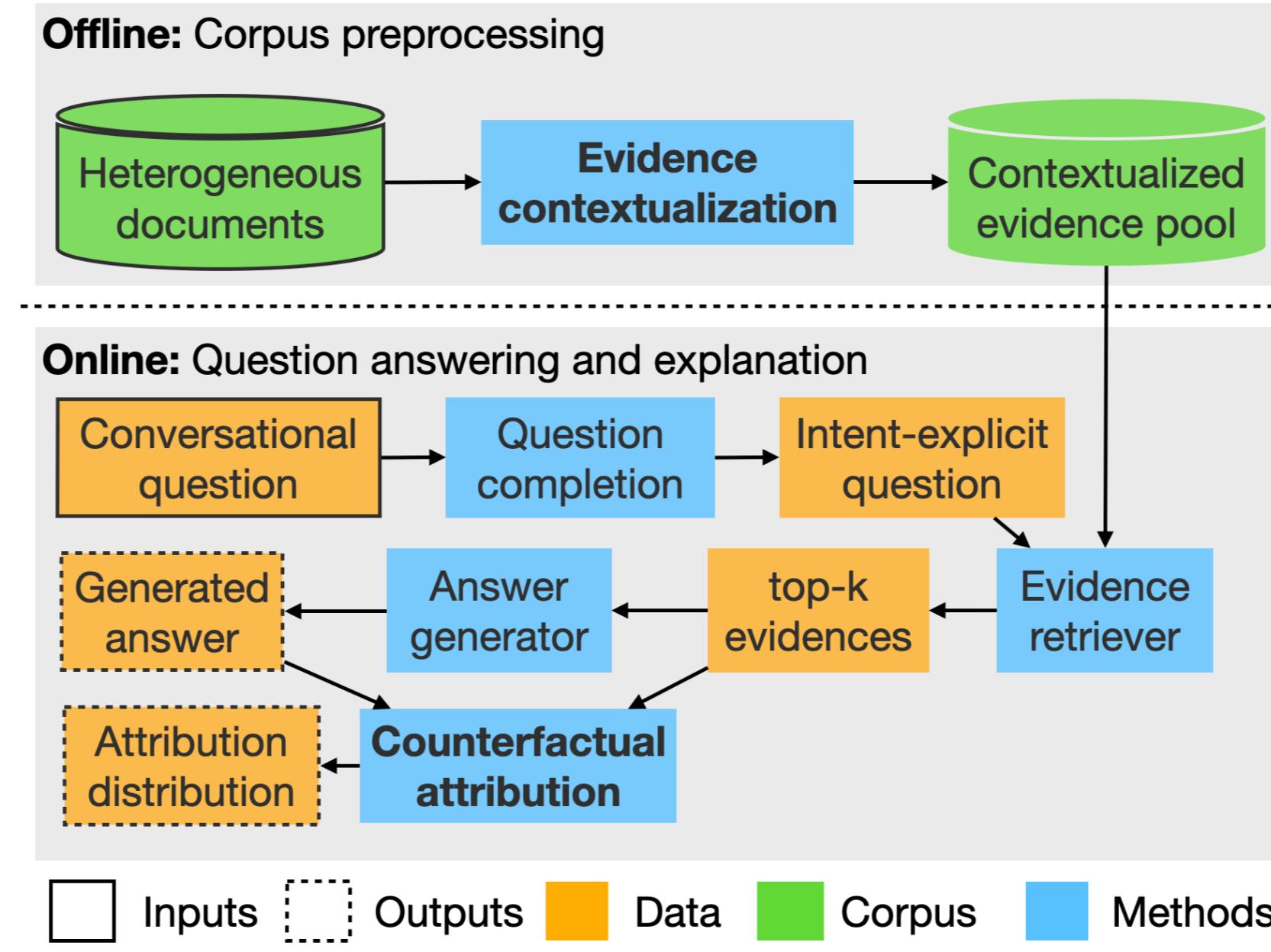
## Idea sketch: Evidence contextualization



## RAGONITE screenshot

The screenshot shows a user interaction with RAGONITE. The user asks, "What was the BIOS and Build versions used for Dell Optiplex 7040 in the OpenXT 9.0 measurement tests?". The system responds with a retrieved passage from a document corpus. Below the main interface, there is a 'Behind-the-scene buttons' section where the user can see the attribution details. It shows a list of evidence with their contribution scores to the final answer. The user also sees follow-up questions and feedback options.

## We propose RAGONITE, and ...

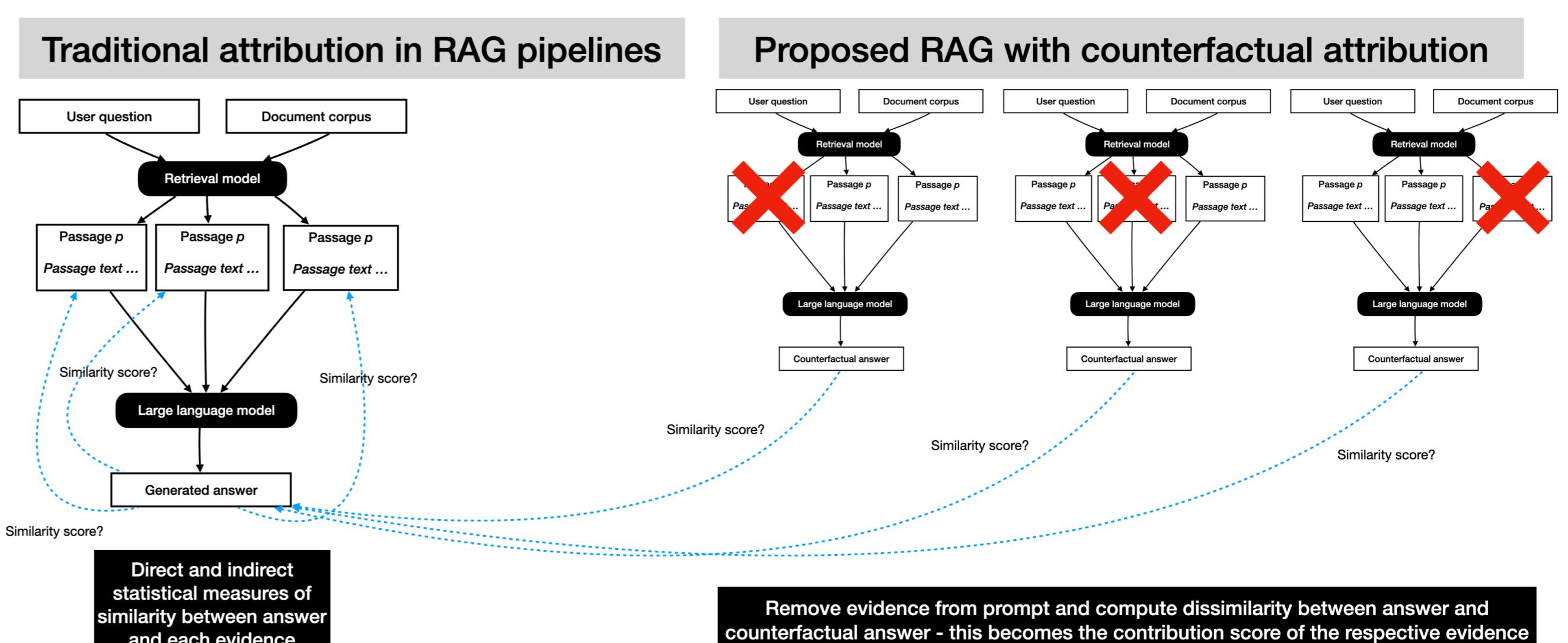


## ... a new benchmark: ConfQuestions

### ConfQuestions is suitable for evaluating

- ✓ Conversational question answering
  - ✓ Complex question answering
  - ✓ RAG on heterogeneous enterprise documents
- ConfQuestions contains**
- ✓ 50 human conversations with 300 questions
  - ✓ Corresponding human-generated completed questions
  - ✓ Gold URL and answers: entities, phrases, lists, passages
  - ✓ Each question in English and German
  - ✓ Metadata: Answer-source, Complexity type
  - ✓ 215 Confluence pages as corpus

## Idea sketch: Counterfactual attribution



## ConfQuestions excerpt

```
"conv_id": "1",
"turns": [
  {
    "turn_id": "1",
    "q_type": "complex",
    "q_en": "What was the BIOS and Build versions used for Dell Optiplex 7040 in the OpenXT 9.0 measurement tests?",
    "q_de": "Welches BIOS und welche Build versionen wurden für den Dell Optiplex 7040 im Zuge der OpenXT 9.0 measurement tests eingesetzt?",
    "completed_q_en": "What was the BIOS and Build versions used for Dell Optiplex 7040 in the OpenXT 9.0 measurement tests?",
    "completed_q_de": "Welches BIOS und welche Build versionen wurden für den Dell Optiplex 7040 im Zuge der OpenXT 9.0 measurement tests eingesetzt?",
    "a_url": "https://openxt.atlassian.net/wiki/spaces/TEST/pages/761823271/OpenXT+9.0+Measurement+Test",
    "a_source": "table",
    "a": "Build version: 6662 and BIOS version: 1.14.0"
  },
  {
    "turn_id": "2",
    "q_type": "complex",
    "q_en": "And what about TPM?",
    "q_de": "Und was ist mit TPM?",
    "completed_q_en": "What was the TPM version used for Dell Optiplex 7040 in the OpenXT 9.0 measurement tests?",
    "completed_q_de": "Welche TPM version wurde für den Dell Optiplex 7040 im Zuge der OpenXT 9.0 measurement tests genutzt?",
    "a_url": [
      "https://openxt.atlassian.net/wiki/spaces/TEST/pages/761823271/OpenXT+9.0+Measurement+Test"
    ],
    "a_source": "table",
    "a": "Version 2.0"
  },
  {
    "turn_id": "3",
    "q_type": "complex",
    "q_en": "Which OpenXT 9.0.2 tests did Dell 9010 pass for the legacy system?",
    "q_de": "Welche OpenXT 9.0.2 Tests hat der Dell 9010 für veraltete (legacy) Systeme bestanden?",
    "completed_q_en": "Which OpenXT 9.0.2 measurement tests did Dell Optiplex 9010 pass for the legacy system?",
    "completed_q_de": "Welche OpenXT 9.0.2 Tests hat der Dell 9010 für veraltete (legacy) Systeme bestanden?",
    "a_url": [
      "https://openxt.atlassian.net/wiki/spaces/TEST/pages/761823271/OpenXT+9.0+Measurement+Test"
    ],
    "a_source": "table",
    "a": "Version 2.0"
  }
]
```

## Evaluation results

Retrieval P@1 for CONTEXTUALIZATION methods						Accuracy for ATTRIBUTION methods		
None	w/ Page title	w/ Preceding heading	w/ Preceding text	w/ Succeeding text	w/ All context	Standard	CF	CF with Clusters
0.398	0.483	0.448	0.453	0.460	0.528 ★	0.772	0.791	0.799 ★
Answer relevance for CONTEXTUALIZATION methods								
None	w/ Page title	w/ Preceding heading	w/ Preceding text	w/ Succeeding text	w/ All context	Standard	CF	CF with Clusters
0.388	0.477	0.435	0.445	0.404	0.529 ★	0.772	0.791	0.799 ★

Take-home message: Contextualizing evidence with surrounding document text, and explaining answers with counterfactual attribution, are vital RAG add-ons



Scan the QR code for all info on project... ☺