# Ask Project Nexodus Docs/Project Aspen
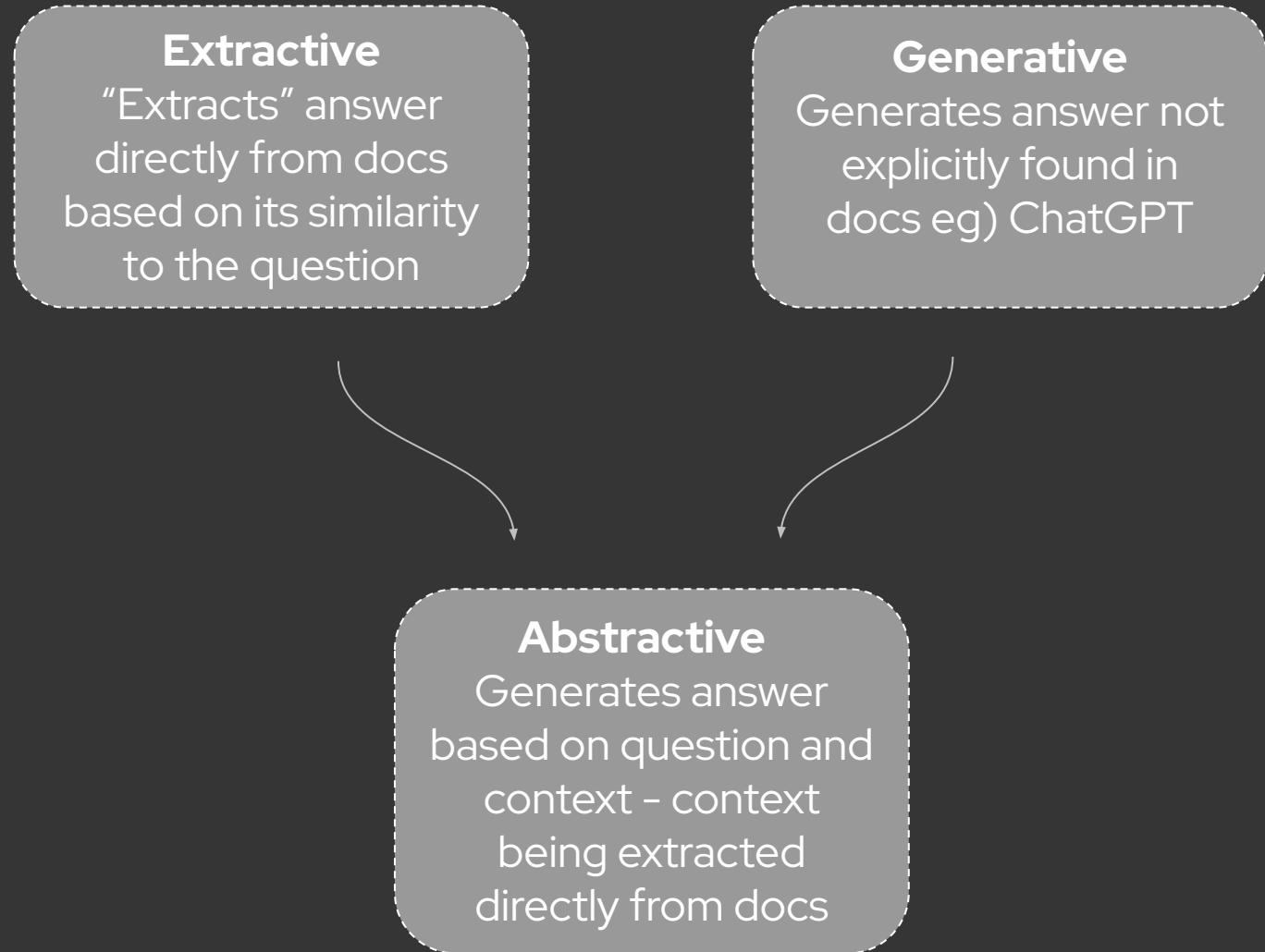
## Intern Final Presentation

Christina Xu

# Ask Project Nexodus Docs: Leveraging LLMs for Documentation Q&A

**Red Hat**

# Main strategies for QA tasks

**Extractive**
"Extracts" answer directly from docs based on its similarity to the question

**Generative**
Generates answer not explicitly found in docs eg) ChatGPT

**Abstractive**
Generates answer based on question and context – context being extracted directly from docs

Red Hat

# What is Fine-tuning

- LLMs are pre-trained on specific domains and tasks such as text generation, question answering, etc. We might want to train the LLM to adapt to our data and task

**Adapter Tuning**

Add more layers to the pre-trained model and train weights only in those additional layers

V0000000

Red Hat

https://dataman-ai.medium.com/fine-tune-a-gpt-lora-e9b72ad4ad3

# What is Fine-tuning?

Traditional approaches are not practical

- LLMs are trained on specific domains and tasks such as text generation, question answering, etc. We might want to train the LLM to adapt to our data and task

**Adapter Tuning**

Ad...
the...
trai... nly
in...
lay...

**Drawbacks**

- Inference latency – the more layers, the longer it takes for the model to generate an answer

5

V0000000

# Fine-tuning with LoRA (Low Rank Approximation)

## Update pretrained weights in the model

V0000000

https://dataman-ai.medium.com/fine-tune-a-gpt-lora-e9b72ad4ad3

# During training...

Decompose ΔW into A and B



Consider a 100 x 100 matrix ΔW. That would mean we would have to train 10,000 parameters. If we decompose it into matrices A and B, which are 100 x 1 and 1 x 100, respectively, we only have 100 parameters to train in each or 200 in total
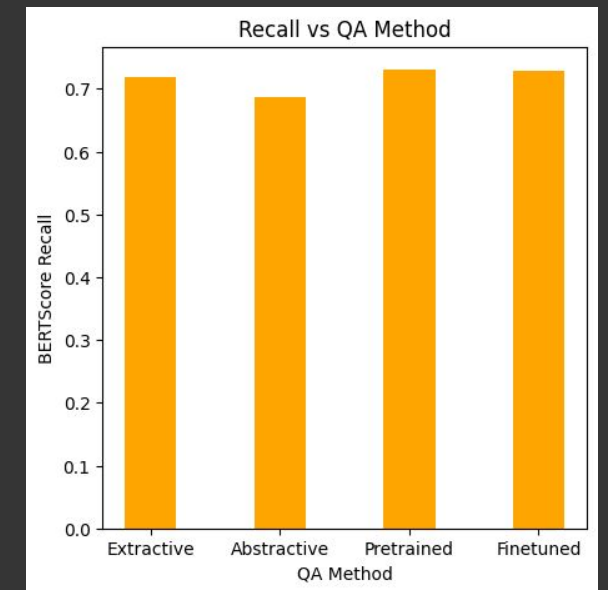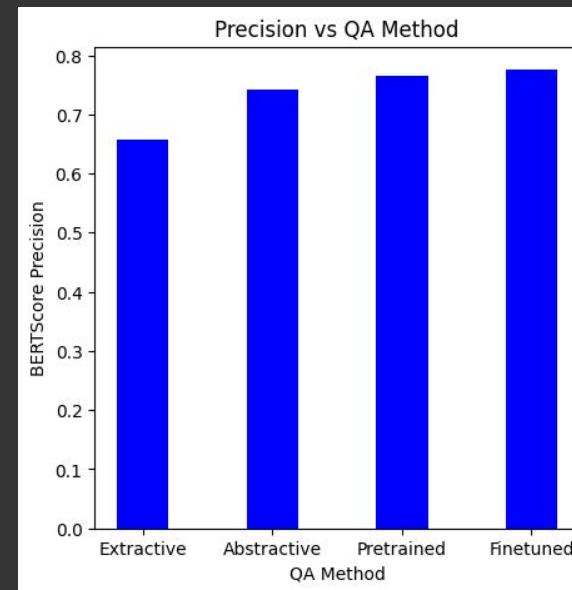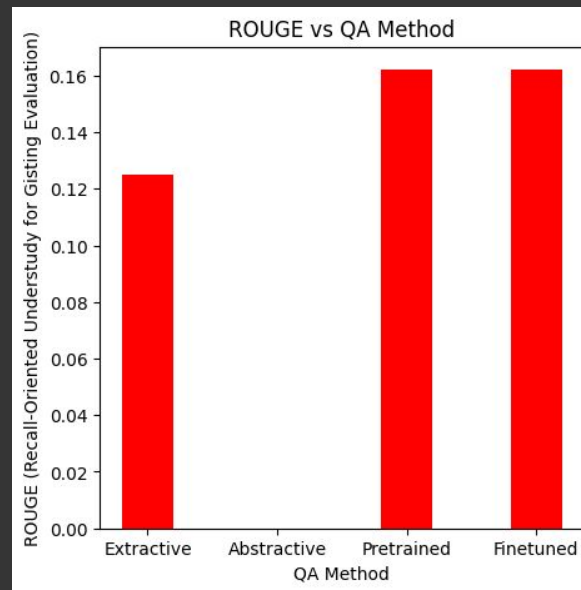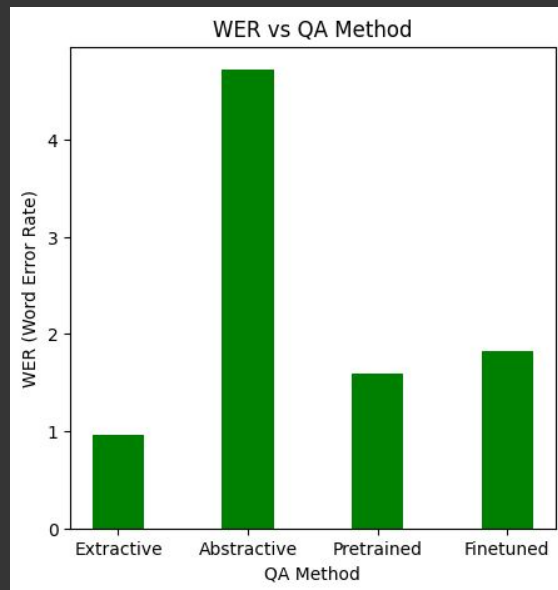
# After training...

## Merge W with ∆W



$$ x \quad \otimes \quad \left( \text{Merged Weights } W + \Delta W \right) \quad \text{output} \Rightarrow \quad y $$

# LLM Strategy Evaluation

Key Idea: human language is difficult to quantitatively evaluate

https://huggingface.co/spaces/evaluate-metric/wer
https://huggingface.co/spaces/evaluate-metric/rouge
https://huggingface.co/spaces/evaluate-metric/bertscore

V0000000

# Project Aspen: Bus Factor

Red Hat

# What is Project Aspen?

Analyzes data from open source projects to empower contributors and participants to make data driven decisions about open source communities and projects.

Red Hat

# Bus Factor

## How high the risk is to a project should the most active people leave?



- Quantifies the amount of contributors a project can afford to lose before it stalls by hypothetically having these people get run over by a bus

- Typically, it is the smallest number of people that make up 50% of contributions

V0000000

https://chaoss.community/kb/metric-bus-factor/

# How do we define "contributions"?

We can analyze bus factors from different perspectives

Commits → Issues → Pull Requests

Red Hat

Top 10 Contributors to the Ansible Repository

**Commits**
- 01012f1b 19.9%
- other 25.8%
- 01000c4d 15.4%
- 01022886 10.4%
- 01000cc2 8.36%
- 01000067 5.72%
- 01000331 5.25%
- 01000e5a 4.23%
- 01001c87 1.75%
- 01000d6a 1.6%
- 010009b6 1.59%

**Issues Created**
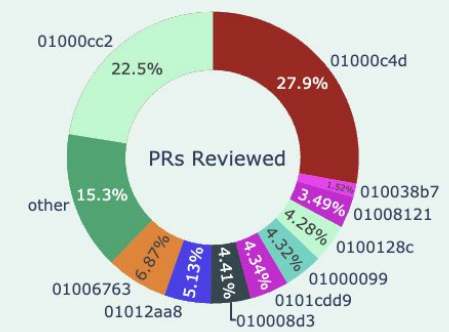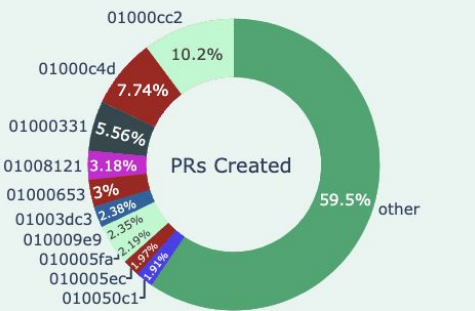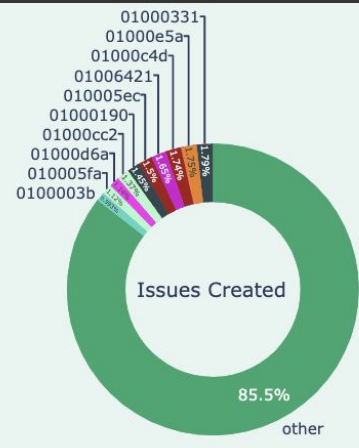- 85.5% other
- 01000331
- 01000e5a
- 01000c4d
- 01006421
- 010005ec
- 01000190
- 01000cc2
- 01000d6a
- 010005fa
- 0100003b

**PRs Created**
- 01000cc2 10.2%
- 01000c4d 7.74%
- 01000331 5.56%
- 01008121 3.18%
- 01000653 3%
- 01003dc3 2.38%
- 010009e9 2.35%
- 010005fa 2.19%
- 010005ec 1.97%
- 010050c1 1.91%
- other 59.5%

**PRs Reviewed**
- 01000cc2 22.5%
- 01000c4d 27.9%
- other 15.3%
- 01006763 6.87%
- 01012aa8 5.13%
- 010008d3 4.41%
- 0101cdd9 4.34%
- 01000099 4.32%
- 0100128c 4.28%
- 01008121 3.49%
- 010038b7 1.52%

**PR Comments**
- 01000067 20.1%
- 01000cc2 12.2%
- 010005fa 4.54%
- 010014b6 3.17%
- 01001c87 2.95%
- 010000ed 2.24%
- 010005ec 1.7%
- 010006cf 1.5%
- 01000ac2 1.45%
- 01000331
- other 48.7%

# Key Insights

- There appears to be a trend in the top 10 contributors across all perspectives

  eg) 01012f1b, 01000c4d, 01000cc2

- The proportion between the top 10 and 'other' contributors for each perspective matches our intuition

V0000000
Red Hat

# Bus factor as a function of time

# Thank you

Let's connect!

Special thanks to Sanjay Arora, James Kunstle, Heidi Dempsey, Jen Stacy, and my fellow research interns

For questions or concerns regarding my projects, feel free to reach out to me via: chrxu@redhat.com

https://github.com/oss-aspen/Rappel

https://github.com/christinaexyou/ask_project_nexodus_docs_(WIP)

https://www.linkedin.com/in/christinaexyou/

https://medium.com/@christinaexyou

Red Hat