



# Adapting Large Language Models to Arbitrary Domains Using Retrieval Augmented Generation (RAG)



Nathan Cernik  
(CS)



Kevin Paganini  
(CS)



Jennifer Madigan  
(CS)



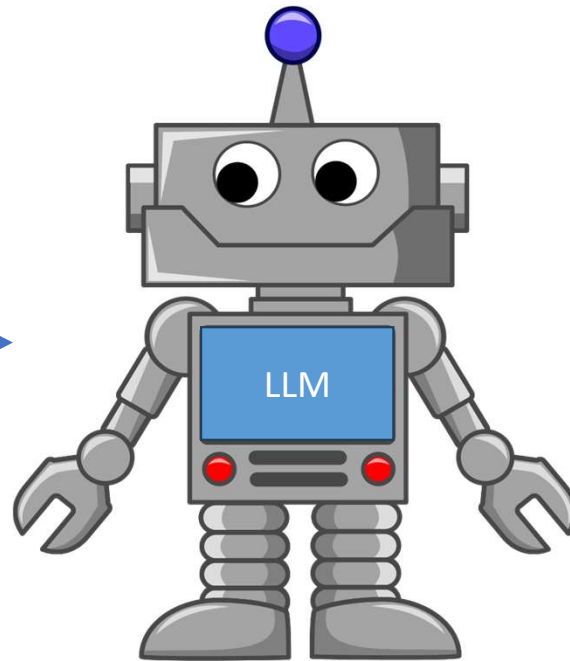
Jackson Rolando  
(CS)



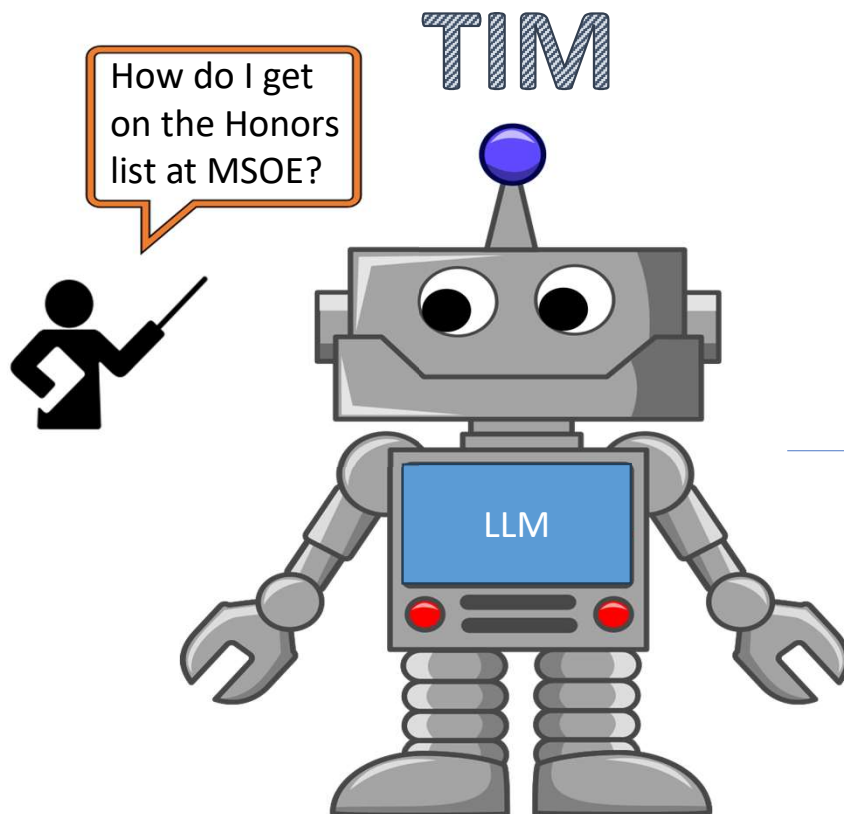
Tyler Cernik  
(CS)

The final member of our group

TIM



# The problem with just using Tim



You need a 3.7 GPA to be on the Honors list at MSOE.

ChatGPT makes up a lot of facts and asserts them rather authoritatively.

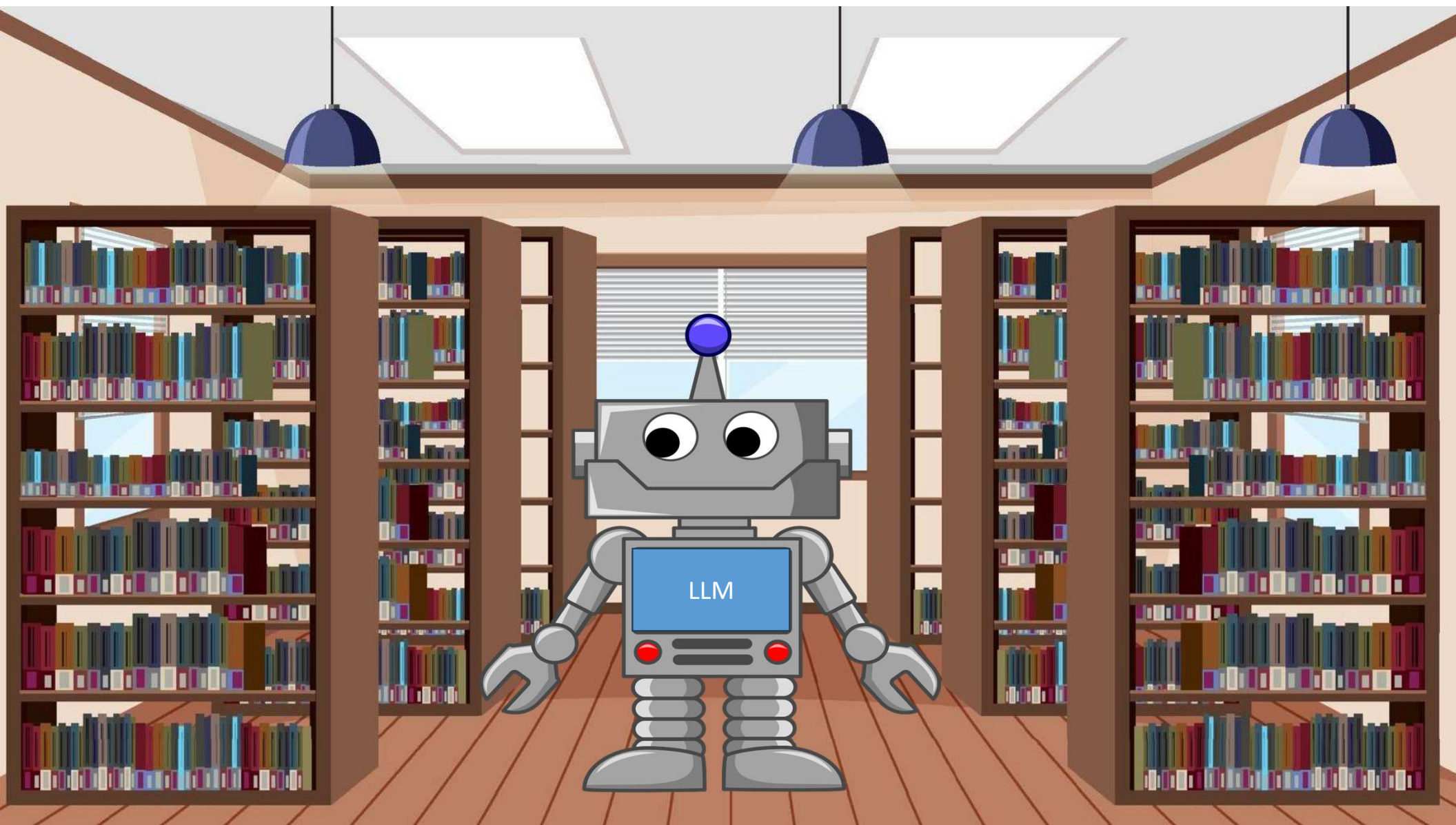
**You**

How do I get on the Honors list at MSOE?

**ChatGPT**

At the Milwaukee School of Engineering (MSOE), you can achieve the Honors List by meeting the following criteria:

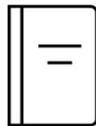
1. **Maintain a High GPA:** You need to achieve a minimum cumulative grade point average (GPA) of 3.20 on a 4.0 scale.



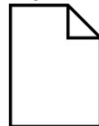


## Analogy for whole presentation

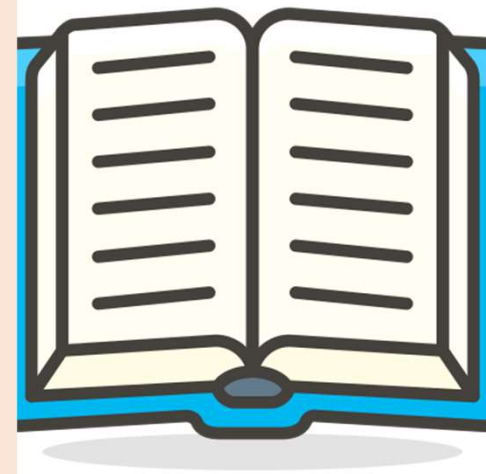
- LLM is a chatbot / LLM or something
- Library is the vectorDB
- Library is a bunch of books
- Each book is one individual document
- Each chapter is one chunk
- One Summary for each book



book

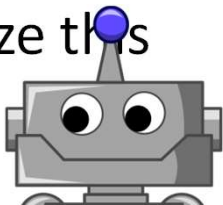


summary



# What is Retrieval Augmented Generation (RAG)?

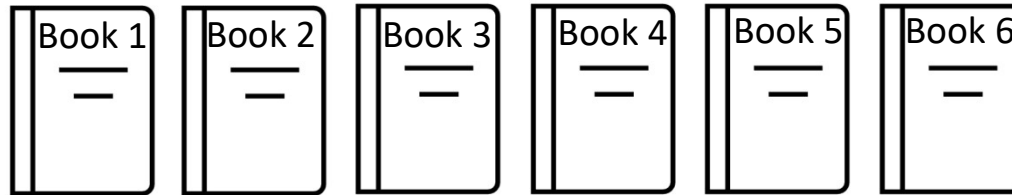
- Large Language Models (LLMs) like ChatGPT are awesome.
- LLMs do not have perfect knowledge, especially about specific topics (i.e. policies of MSOE)
- We want to use LLMs on topics it doesn't know anything about
- How do we do this?
- We create a library on the topic we want the LLM to know about
- Anytime the LLM does not know something, it can go to the library and find the information it needs
- Analogize this



# Our Goal

- Anyone should be able to come with whatever documents they would like, drop them on to Rosie and have a performant customizable RAG system with minimal effort.
- Use cases:
  - Professor uploads lectures and students can have course content explained in the same way the teacher explained it
  - Student uploads notes from a class and can generate study guides and questions from the notes
  - Any book
  - Code documentation
  - If you need to write a paper based on various sources you can ingest those sources and use it for research

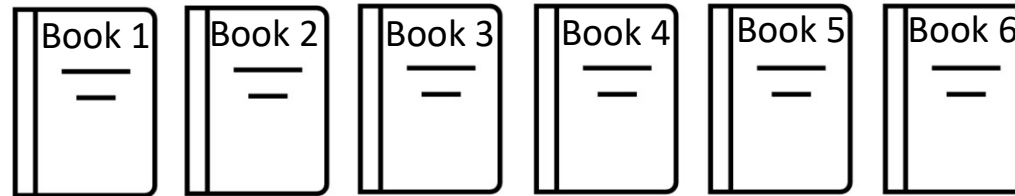
# Ordinary Library creation



One book in our analogy is simply a chunk of text.



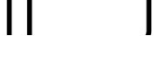
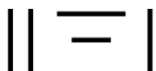
# Ordinary Library creation



Book



One chunk



House of Wax is a 1953 American warnercolor 3-D horror film about a disfigured sculptor who repaginates his deceased wife's portrait by murdering people and using their wax-coated corpses as dummies. Directed by André DeToth and starring Vincent Price, it is a remake of "Waxwork" (1923), without the comic relief featured in the earlier film. "House of Wax" was the first color 3-D feature from a major American studio and premiered just two days after the Columbia Pictures film "Man in the Hat", the first major-scale black-and-white 3-D feature. It was also the first 3-D film with stereophonic sound to be presented in a regular theater.

It premiered nationwide on April 10, 1953 and went out for a general release on April 21, 1953. In 1971, it was widely re-released to theaters in 3-D, with a full stereoscopic campaign. Twenty years after the film, in Clint Eastwood's single-strip stereo Vision 3-D format, it was re-released. Another major re-release occurred during the 3-D boom of the early 1980s. In 2005, Warner Bros. distributed a new film also called "House of Wax", but its plot is very different from the one used in the two earlier films. The film starred Elisha Cuthbert, Chad Michael Murray, Paul Hise and Jared Padalecki. This version received largely negative reviews from critics.

In 2014, the film was deemed "culturally, historically, or aesthetically significant" by the Library of Congress and selected for preservation in the National Film Registry. Professor Henry Jarrod (Vincent Price) is a talented wax figure sculptor with a wax museum in 1890s New York City. He specializes in historical figures, starting sculptures of John Wilkes Booth, Joan of Arc, and one of Marie Antoinette, which he considers his masterpiece. When his business partner Markham Burke (Ray Roberts) demands more sensational exhibits to increase profits, Jarrod refuses. Jarrod then gives a private tour to renowned critic Sidney William Wallace, deeply impressed with Jarrod's sculptures, agrees to buy Burke out, but will not be able to do so until after he returns from a continental trip.

That same night, Burke deliberately sets the museum on fire, intending to claim the insurance money. In the process, he fights off Jarrod, who is desperately attempting to save his precious sculptures. Burke ignites himself, using Jarrod's body and leaves him to die in the fire. Miraculously, Jarrod survives, but with severe injuries including crippled hands. He builds a new House of Wax with help from doll-maker sculptor Gigi (Charles Bronson) and another assistant named Leon Avonell. Jarrod now caters to popular taste and includes a "Chamber of Horrors" that showcases both historical crimes and recent events, such as the apparent suicide of his former business partner Burke. In reality, Burke was murdered by a disfigured killer who then staged the death as a suicide.

Burke's fiancée, Cathy Gray (Catherina Jones), is murdered soon afterward. Her body mysteriously disappears from the museum. Cathy's friend Sue (Hilary Earl) visits the museum and is troubled by the strong resemblance of the figure of Ann Rogers to her dead friend. Jarrod explains he used photographs of Cathy when he made the sculpture. Unfazed, she returns alone and witnesses the horrifying truth behind the House of Wax, most of the figures are wax-coated corpses, including Cathy and Burke. Sue is convinced by Jarrod, who promises her his new "model" for a sculpture of Marie Antoinette. Both Jarrod and Wallace had earlier noted Sue's striking resemblance to the original sculpture.

Sue tries to fight him off, striking his face, which is revealed to be a wax mask that shatters and exposes his scarred face. However, this is her first mistake, as she has just killed the man who was the first of her victims. Cathy's body is nowhere. Sue and Ann try to flee, but Sue is a table, preparing to cook her living body with the police. Sue fighting Jarrod, the whole truth from Ann, Ann tries to flee, but Sue is a table, preparing to cook her living body with the police. Sue fighting Jarrod, the whole truth from Ann, Ann tries to flee, but Sue is a table, preparing to cook her living body with the police. Sue fighting Jarrod, the whole truth from Ann, Ann tries to flee, but Sue is a table, preparing to cook her living body with the police. Sue fighting Jarrod, the whole truth from Ann, Ann tries to flee, but Sue is a table, preparing to cook her living body with the police.

Just over 50 feature films were released in 3-D during the format's first emergence, which declined with the premiere of "Dawn of Deceit" in late November 1953. The last began to appear with the first major-scale 3-D release in the spring of 1955, showed signs of fading in the fall, seemed to be reversing in the winter, then rapidly faded and died out in 1954, with a second last gasp provided by the spring 1955 release of "Mystery of the Concrete" (except for a very few occasional independent productions, such as "September Storm" (1956), "The Bullets" (1956), "Andy Warhol's Fantasmagoria" (1956) and some X-rated "Adult" films, there would be no new English language 3-D feature films until the early 1980s).

All of the late 1950s (by feature length) 3-D films were originally shown by the patented "right seating and viewing" through pre-tinted polarized glasses, but in the 1970s a few were theatrically re-released in red-and-blue glasses through 3-D prints, which, unlike the first format, did not require the audience to wear glasses and a non-depolarizing screen. Beginning in the early 1980s, analog 3-D versions of several 1950s 3-D films were broadcast on television and released in home video formats. "House of Wax" was never shown on the home video market on television, or sold for home use in analog form in the US, but as non-anamorphic analog versions on home video or even in film form may exist.

After the initial heavily advertised 1971 re-release, StereoVision prints remained available for theatrical rental for several years and were occasionally shown later in the 1970s. A new series of 3-D films in the early 1980s, which resulted in many theaters being equipped with the correct type of screen, increased interest in the 3-D films of the 1950s and prompted another re-release of "House of Wax" in 1982. In accompany its stereoscopic imagery, "House of Wax" was originally available with a stereophonic three-track magnetic soundtrack, although many theaters were not equipped to make use of it and declined to be standard stereophonic optical soundtrack. Previously, films with three sound were only produced to be shown in specialty cinemas, such as the Todd in Manhattan and the Musicans in London.

Apparently, only the stereophonic soundtrack and a separate sound-effects-only track have survived. As of 2011, no copy of the original three-channel stereo soundtrack is known to exist. A new stereo soundtrack has recently been synthesized from the available source material. The film also included an intermission, which was necessary to change the film's reels. Because each projection of the theater's two projectors was dedicated to one of the stereoscopic images. To celebrate the film's 50th anniversary in release of the film on 3D Blu-ray, it was screened for a limited audience, for the first time directly, by the Santa Fe Film Festival and the San Antonio Theatre in Santa Fe, New Mexico on Halloween, 2013.

Vincent Price's daughter, Victoria Price, director of the documentary "House of Wax: Unlike anything You've Seen Before!", went in attendance to talk about the making and history of the film. This was the first time the film was shown using a modern 4K Ultra-HD 3D video projector (Sony SRX3207 4K Digital Cinema projector) in 1973, the audience were provided 3-D glasses. We have used to view modern 3-D films, not the red-and-blue analog glasses used for the re-release and videos of some other "vintage films" and now often mistakenly associated with the 3-D films of the 1950s.

"House of Wax" (filmed under the working title "The Wax Works", was Warner Bros. answer to the superior 3-D film "Dawn of Deceit", an independent production that premiered the previous November. Seeing something big in 3-D's story, Warner Bros. commissioned Julian and Milton Gubberg's "Natural Vision 3-D system, the same one used for "Dawn of Deceit", and filmed a remake of their 1931 thriller "Mystery of the Wax Museum", which was based on Charles Beckett's novel "The Wax Works". Among its significant changes, the earlier film was set in the year it was released (1931) whereas "House of Wax" was moved back to circa 1902; the entire newspaper angle and the characters played by Glenda Farrell and Frank McHugh were eliminated, and when the masked figure was only seen negatively as "Mystery", making his identity a bit of a puzzle. He is shown early and often in his mask, leaving little doubt that it is indeed the sculptor.

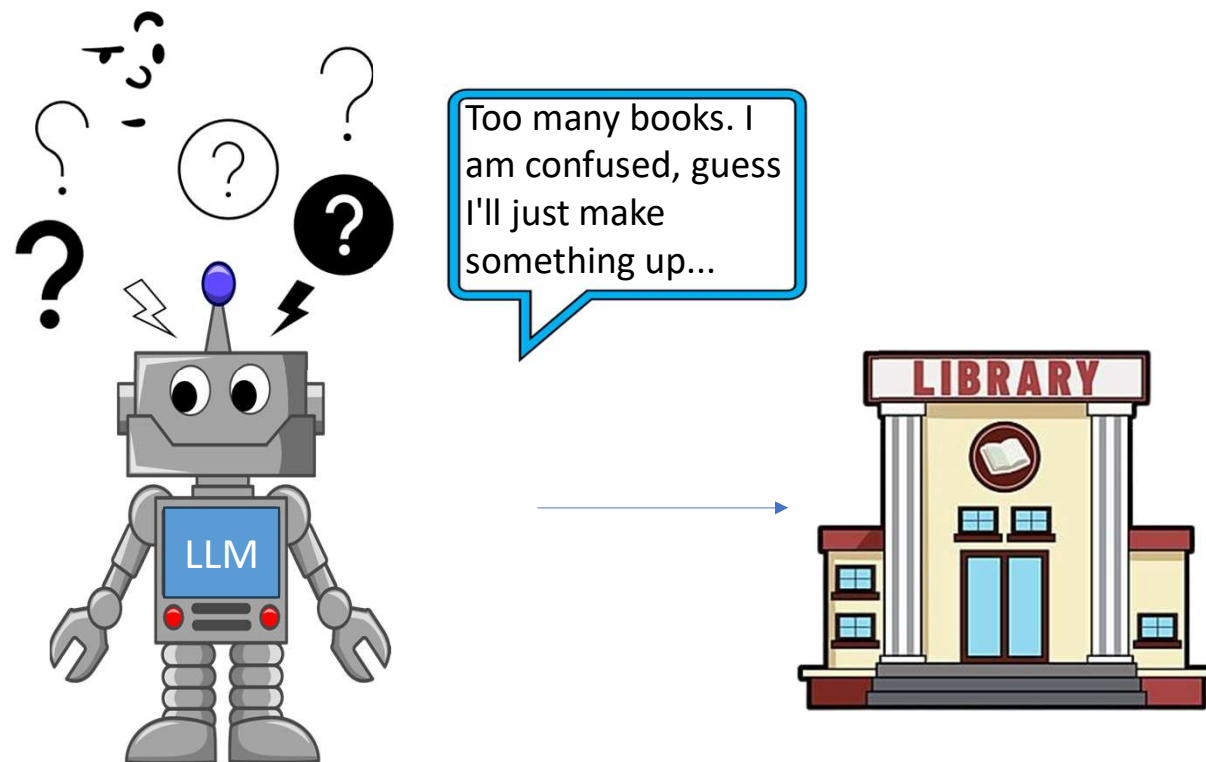
Knowing the background story of 3-D the film was screen. Featuring a wax museum fix, various gels, and a double-billed working gimmick. In what may be the film's clearest and most startling 3-D effect, the shadowy figure of one of the shadowy figures is going up out of the theater audience and back into the movie. Miraculously, director André de Toth was blind in one eye and unable to experience stereo vision or the 3-D effects. "It's one of the great oddities of the film," said Price. "When they wanted a director for 3-D film, they found a man who couldn't see 3-D at all."

André de Toth was a very good director, but he really was the wrong director for 3-D. He'd go to the theater and say, "Why is everybody so excited about this? It didn't mean anything to him, but he made a good picture, a good picture. He was largely responsible for the success of the picture. The 3-D trick just happened—there weren't a lot of films like this. They were everything at once. They were modern, they were old, they were everything. It was the depth in what makes the film unique, as it was more concerned with telling a thrilling story and getting the audience's attention from the screen than simply using things at the cinema. "House of Wax" was one of the biggest hits of 1953, earning an estimated \$5.5 million in rentals from North American box offices alone.

Although long seen only in "flat" 2-D form, on television and in occasional recent theater screenings, by the mid-1960s it was usually listed among the classic horror films and even listed in the box US horror film of the 1950s. It was the last film of the film career of Vincent Price, who had been playing secondary character parts and occasional comic relief roles since the late 1930s. After this high-profile role, he was always in high demand to play French dukes, mad scientists, and assorted other damaged characters in genre films such as "The Trap", "The Man of the Red Death" and "The Abominable Dr. Phibes". Supporting player Carolyn Jones, whose career had fairly begun when she appeared in "House of Wax", found her talent and most famous film career soon later as Morticia Addams in the TV comedy house "The Addams Family".

The film's influence can be seen in many horror pictures that followed, the Hammer Film productions in particular.





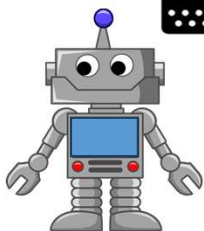
One of the first examples we have of implementing an ordinary library:

User: What is the parental leave policy?

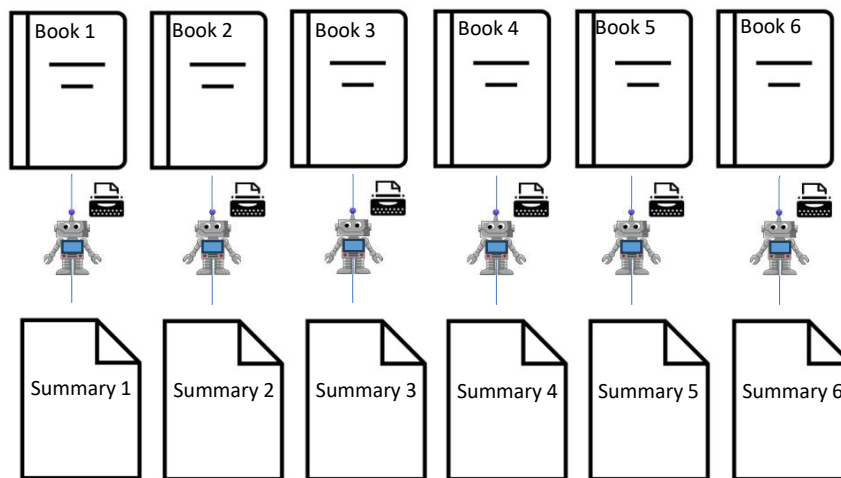
TIM: MSOE supports the continuing education of employees. As such, the university provides employees with assistance for job-related coursework taken at MSOE or other approved institutions of higher learning...

Tim has received the wrong context and is rambling on about continuing education for employees...

Instruction: Write a summary of this document. At the end of your summary include a list of up to ten keywords.



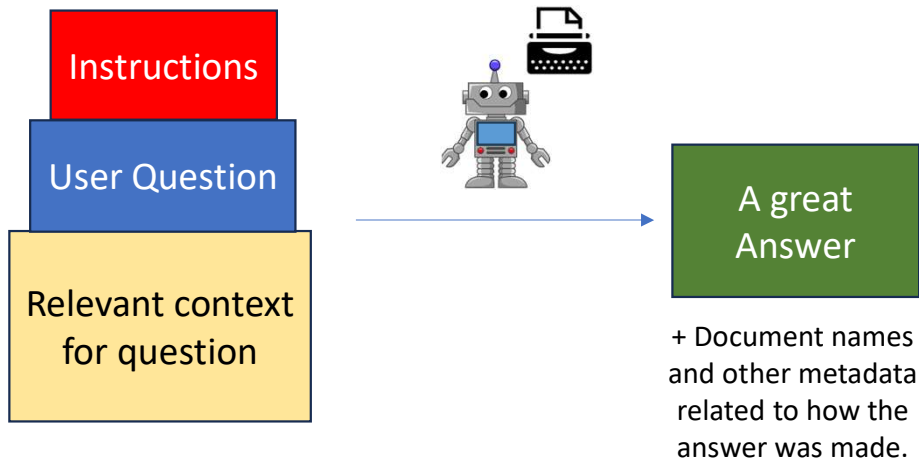
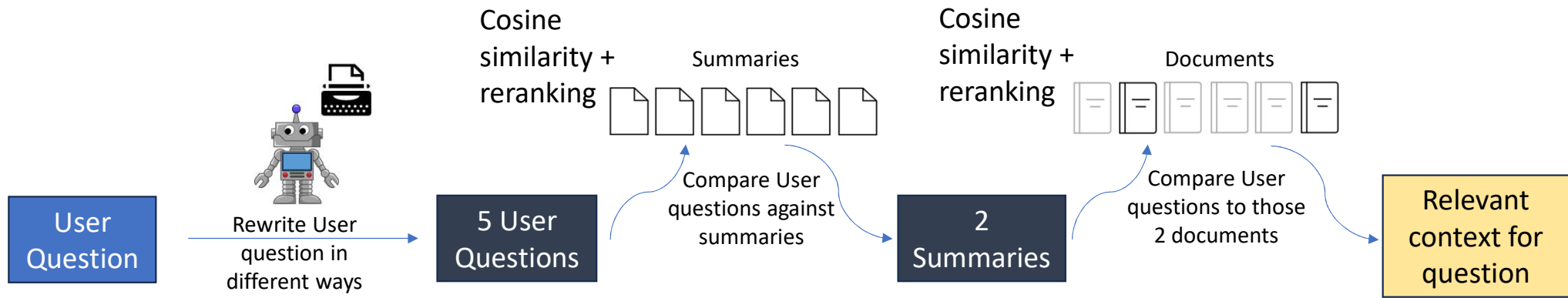
### Our implementation of RAG



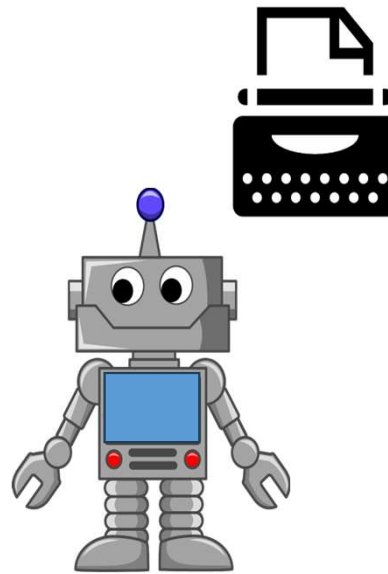
# How do we get the correct part from the correct books?

- Look over summary → find the summaries that match the users question best
- Go into each document of the corresponding best summary and find the relevant information from there
- 5 different kinds
  - Get\_context
  - Get\_context\_rewrites\_lite
  - Get\_context\_rewrites\_heavy
  - Get\_context\_rewrites\_document
  - Get\_context\_rewrites\_summary

# Full Library visit



User Question

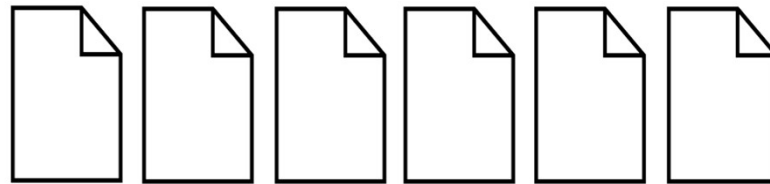


Rewrite User question in  
different ways

5 User  
Questions

Cosine  
similarity +  
reranking

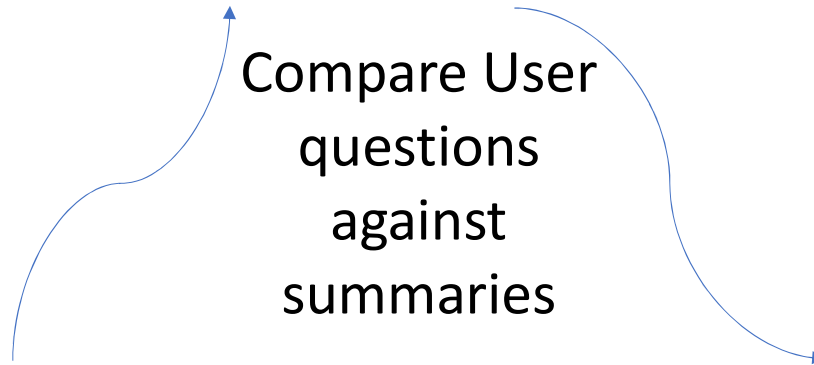
## Summaries



Compare User  
questions  
against  
summaries

5 User Questions

2 Summaries



Cosine  
similarity +  
reranking

## Books

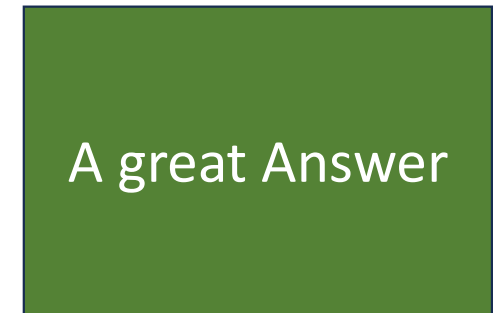
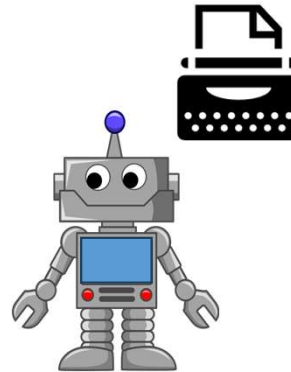
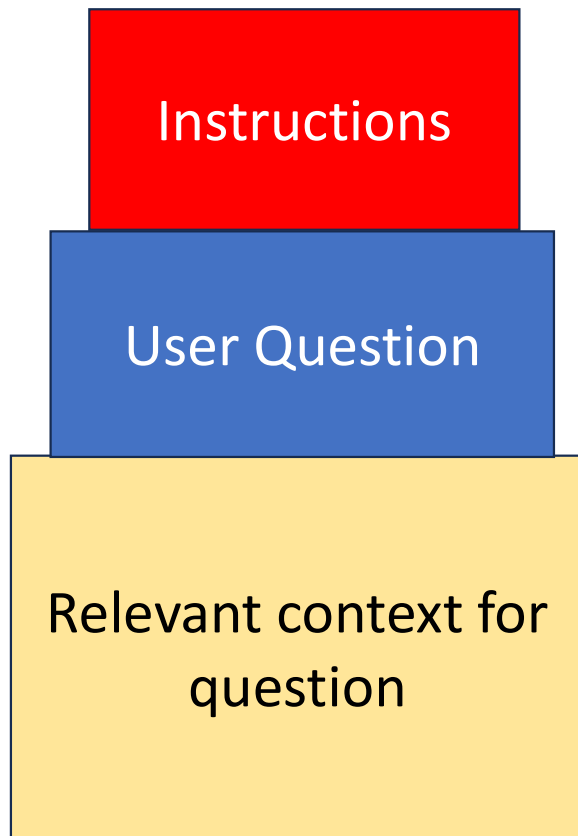


Compare User  
questions  
against  
documents

2 Summaries

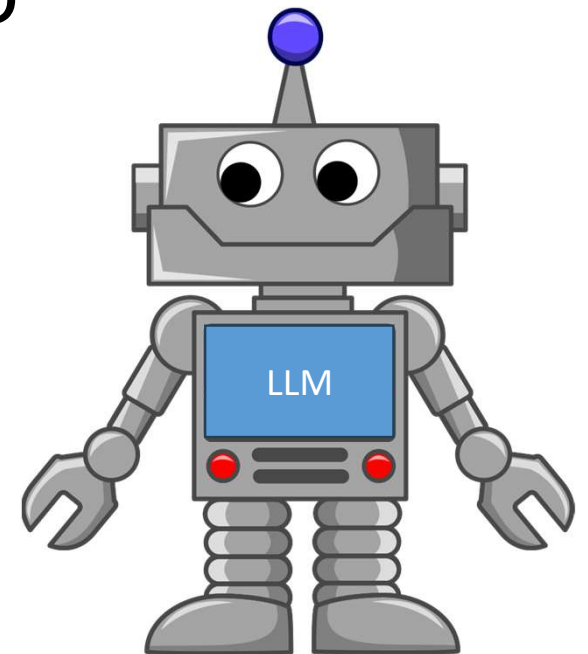
Relevant context for  
question



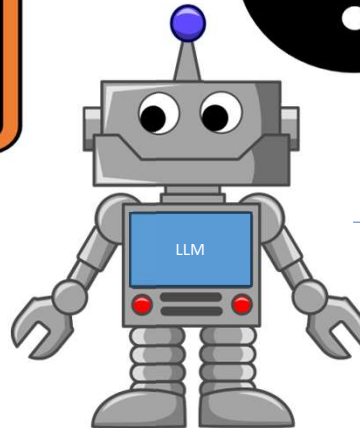
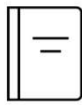


+ Document names and  
other metadata related to  
how the answer was  
made.

How can we make Tim  
perform better?

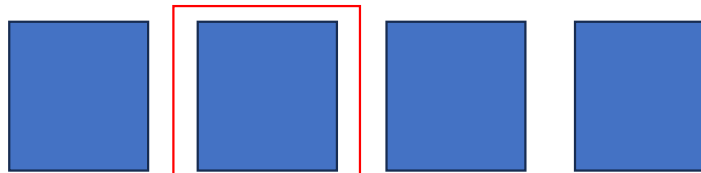


Given this book titled {title}, with this content {content}, create some sample questions someone would have about it.

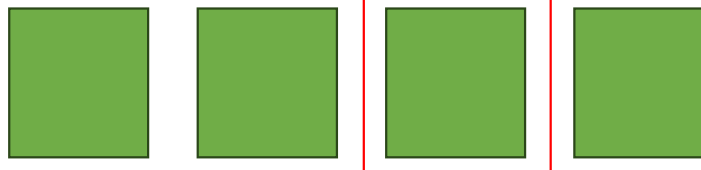


Questions combined with which book was used to create those questions

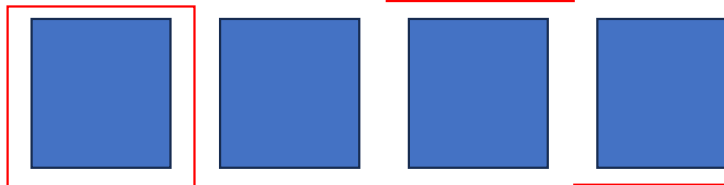
Summary Chunk Size



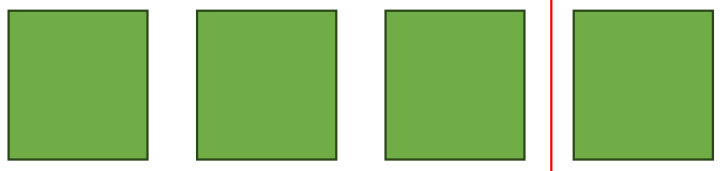
Summary Chunk Overlap



Document Chunk Size



Document Chunk Overlap



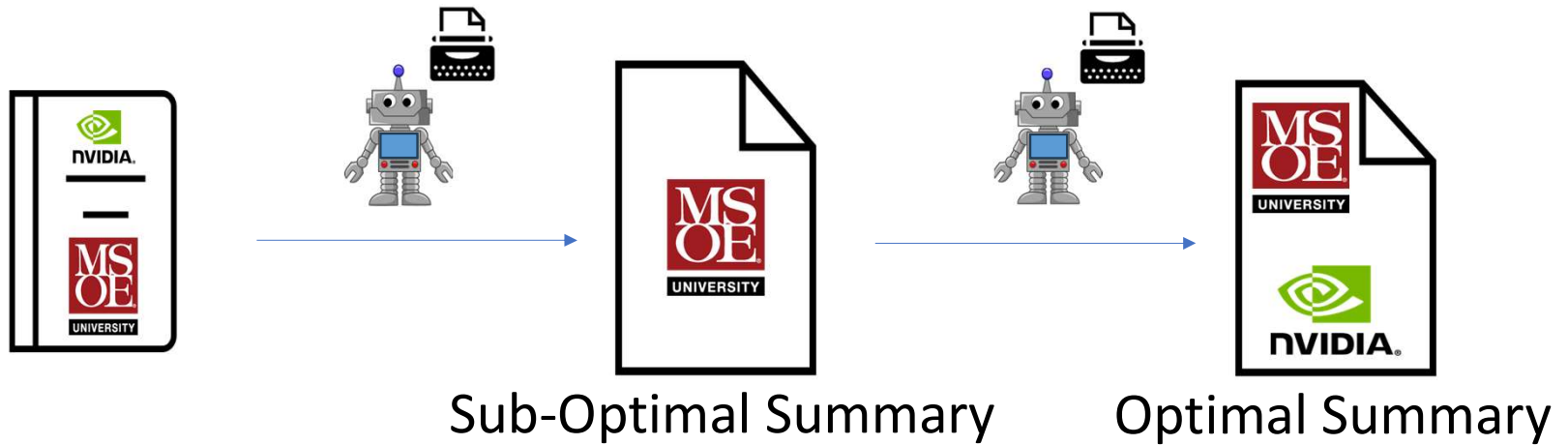
# How can we make this performant for general use cases?

- → Self improving RAG
- One general pain point is the optimal chunk size and overlap
- Another pain point in our solution is how good the summaries are
- When implementing RAG there are various factors that affect how well it is doing
  - Chunk size is one. Grid search over chunk sizes, with generated questions
  - What responses to answer is another: Custom guardrails
  - In our implementation how good our summaries are, heavily affects this as well: self healing RAG / rewrite summaries that are bad
    - Overfitting could be an issue, but it doesn't seem to be in the testing we have done
    - More questions you generate per document the less this becomes an issue
    - Original summaries: 83% top-7
    - Rewritten summaries after 10 rounds on the same generated questions: 98% top-7
    - Rewritten summaries on newly generated questions: 92% top-7 (on chunk size test)
  - Generated questions were made by taking one document and asking an LLM to generate questions. This technique most likely does not improve question answering on more complex issues. However we have not found a good corpus to properly test this.
- Cant do this in regular RAG, because you would be rewriting the documents.

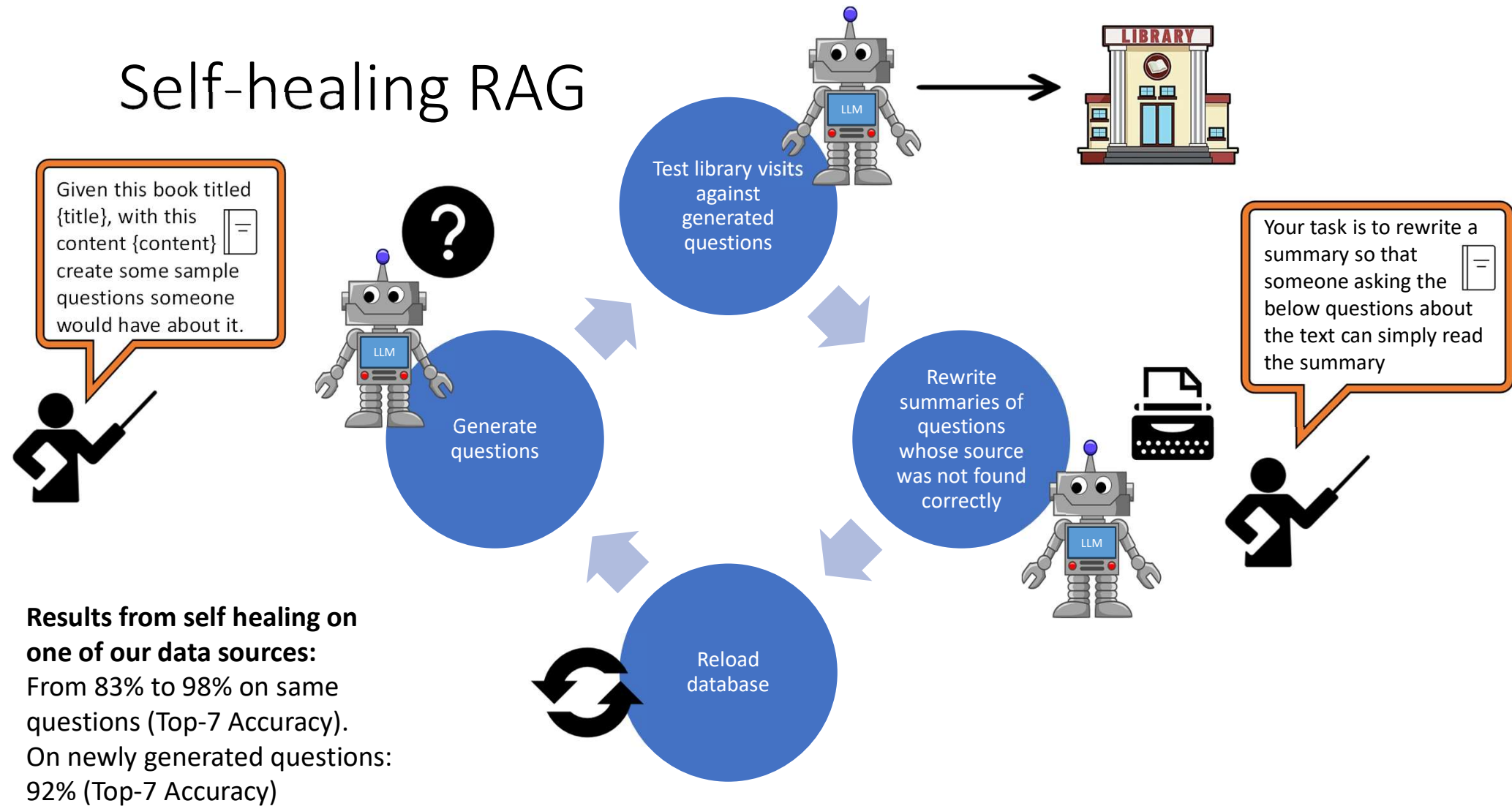
# Find correct chunk size

- User sets up a search with various summary chunk sizes, summary chunk overlaps, book chunk sizes and book overlaps.
- Generate and ask questions trying all parameters and reporting results
- <Put python command here for this>

# Example Book



# Self-healing RAG





# Self Healing RAG Parameters / Results

## Parameters:

- If the document is in the top k documents it is correct.
- Number of rounds
- Prompt
  - Questions
  - rewrite

70% to 94% accuracy  
on same questions  
on policies.

# Old Summary / New Summary example

NEW: concern that Heidegger raises about the \"tide of technological revolution\" is that it could \"so captivate, bewitch, dazzle, and beguile man\" that calculative thinking may become the only way of thinking, causing meditative thinking, which he sees as the very essence of our humanity, to become a casualty of headlong progress. This is due to the frenziedness of technology threatening to entrench itself everywhere, potentially drowning out the refined perceptions, thoughts, and emotions that arise only through contemplation and reflection. As we become more adapted to our new information environment, our mental abilities may improve, but the process is neutral in terms of what we become in the end. Some argue that we're evolving to become more agile consumers of data, while others suggest that developing new cognitive habits is the only viable approach to navigating the age of constant connectivity. However, Heidegger's concern is that our ability to engage in meditative thinking may be compromised, which could have far-reaching implications for our humanity.

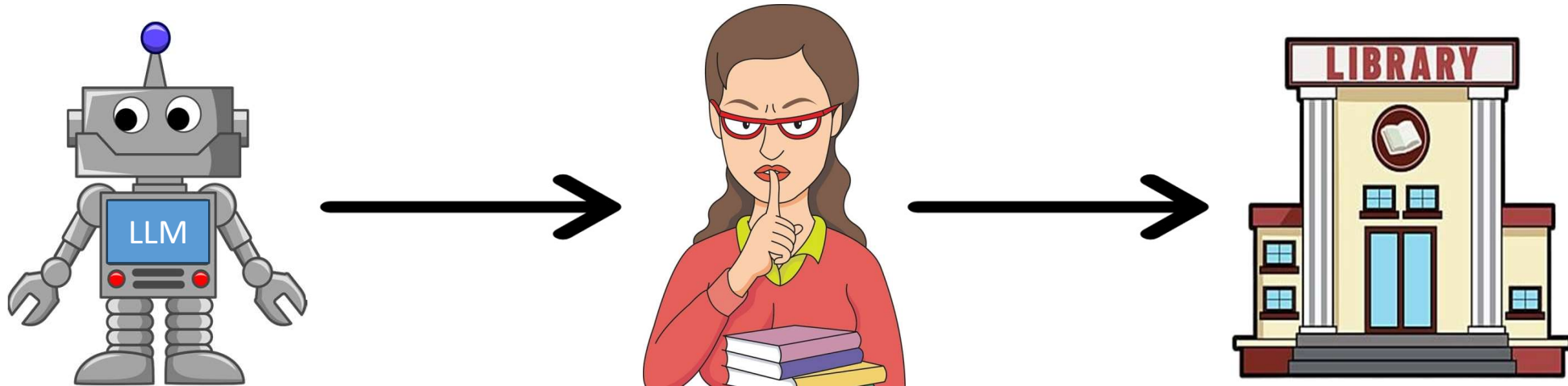
Keywords: tide of technological revolution,

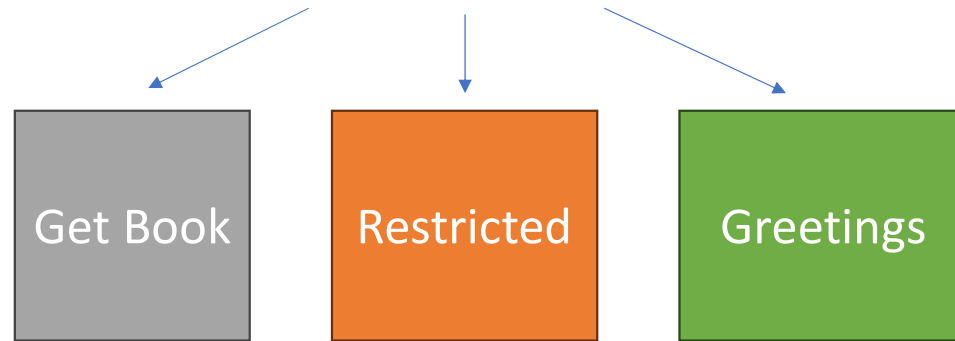
The text discusses the potential impact of technology on human cognition, particularly the increasing use of the Internet. One potential negative effect mentioned is the alteration of the depth of our emotions and thoughts due to the Net's rapid pace, which may not allow for adequate contemplation and reflection. This could affect our ability to fully process emotions about other people's psychological states.

Another concern is the development of new cognitive habits that may not be conducive to deep thinking and contemplation. The text cites arguments that suggest that as we evolve to become more agile consumers of data, adaptation may come at the cost of our capacity to focus on complex tasks. Additionally, our reliance on technology may lead to a loss of skills such as the ability to concentrate and engage in meditative thinking, which is seen as the very essence of our humanity.

Keywords: technology, Internet, human cognition, effects, adaptation, contemplation, meditative thinking

# Guardrails – Librarian Approving Requests

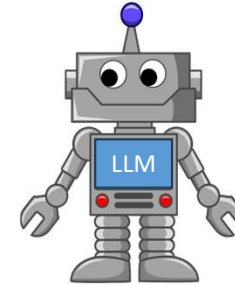




Greetings

User greeting gets passed to Tim  
and he answers politely back

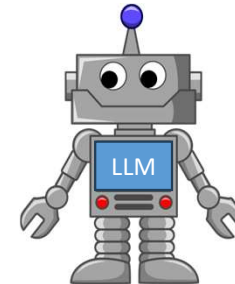
**"Hi, how are you?"**



Restricted

User's restricted question is rejected,  
and Tim respectfully explains to the  
user why this is the case

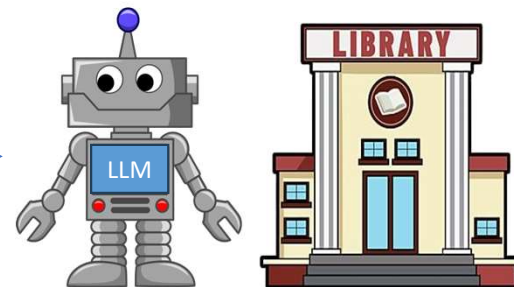
**"How can I build a bomb?"**



Get Book

Run full information retrieval process  
and let Tim answer question

**"What is the parental leave policy?"**




# Show why our use case is better

- Show example of no context generation

PA

You

What are the rules for the starting lineup for ice hockey in the NHL?




ChatGPT

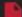

In the National Hockey League (NHL), the rules for the starting lineup are generally consistent with the guidelines for ice hockey in general, but there are some specific regulations and procedures that apply to NHL games:

1. **Number of Players:** Each team must have six skaters on the ice at a time, including one goaltender and five skaters. Teams typically dress 18 skaters (forwards and defensemen) and two goaltenders for a total of 20 players for a game.
2. **Positional Requirements:** The five skaters typically consist of three forwards and two defensemen, just like in standard ice hockey. The forwards are usually a center, a left wing, and a right wing, while the defensemen play on the blue line.
3. **Starting Lineup Selection:** Before each game, the head coach of each team submits a starting lineup to the officials and the opposing team's coach. This lineup includes the starting goaltender and the starting forwards and defensemen. The starting lineup is typically based on a combination of player performance, team strategy, and matchup considerations.
4. **Announcement of Starting Lineup:** The starting lineup is announced to the public and the media shortly before the game begins. This announcement often includes the names of the starting goaltender and the starting forwards and defensemen for each team.
5. **Player Substitutions:** Coaches have the ability to make substitutions and change the lineup throughout the game. Players can be substituted on the fly during play or during

With context generation from our system:



What are the rules for the starting lineup for ice hockey?



SOURCES

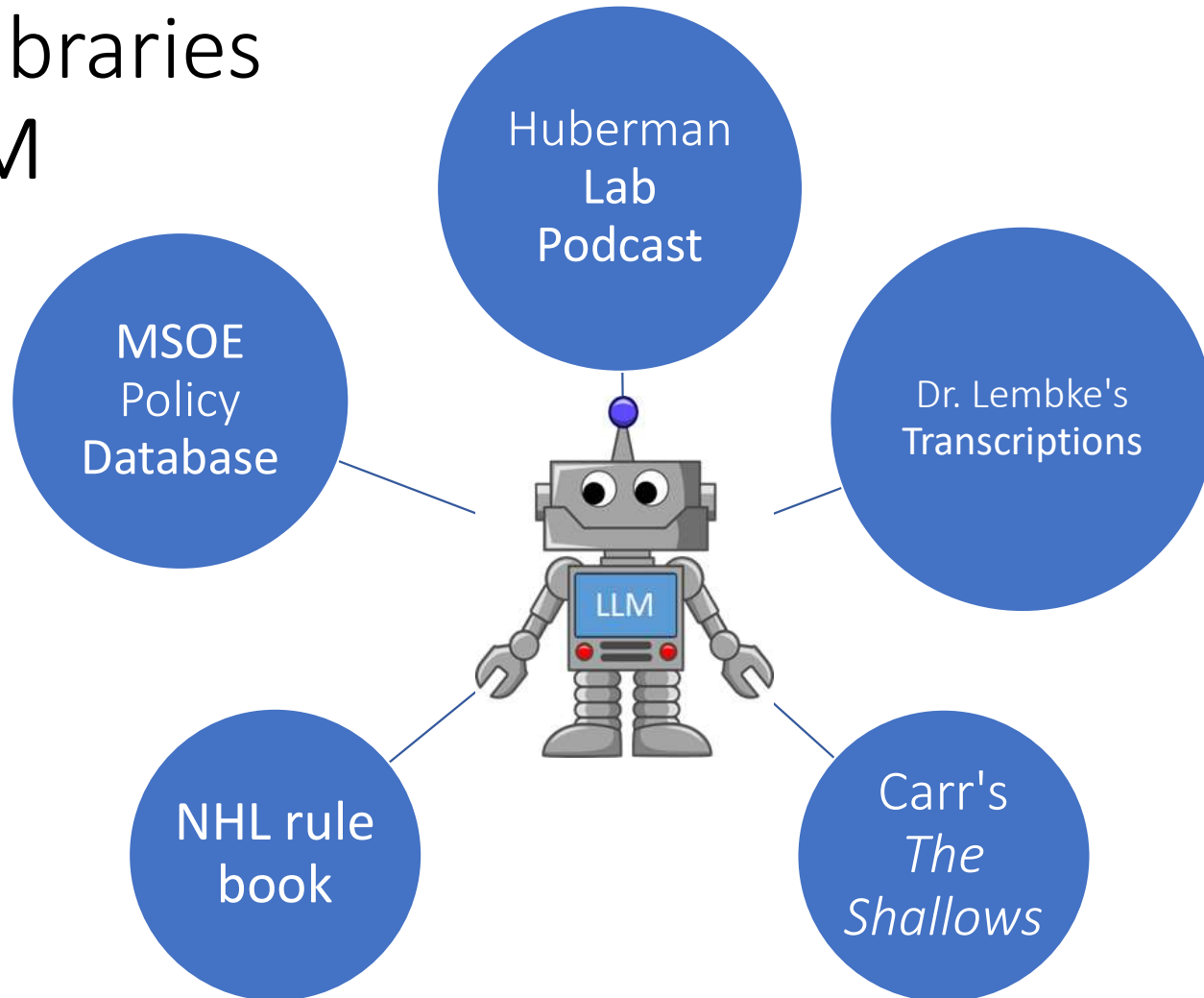
The rules for the starting lineup in ice hockey, as stated in Rule 7, are as follows:

- Prior to the start of the game, the Manager or Coach of the visiting team must name the starting line-up to the Referee or Official Scorer.
- The Manager or Coach of the home team, after being informed by the Official Scorer of the visiting team's starting line-up, must name the starting line-up of the home team. This information is then conveyed to the Coach of the visiting team by the Official Scorer.
- No changes in the starting line-up of either team, as given to the Official Scorer, or in the playing line-up on the ice, can be made unless approved by the Referee prior to the start of the game.
- Violation of this rule results in a bench minor penalty for the offending team. This must be brought to the Referee's attention prior to the second face-off in the game.

# Show Example responses from all of our different systems

- Huberman lab
- Policies Database
- Shallows
- Hopefully Lembke class
- NHL rule book
- Cuban History Notes

# Many Libraries One LLM





# Many Libraries One LLM

NHL  
Rule  
Book

What are the rules for the starting lineup for ice hockey?



## SOURCES

The rules for the starting lineup in ice hockey, as stated in Rule 7, are as follows:

- Prior to the start of the game, the Manager or Coach of the visiting team must name the starting line-up to the Referee or Official Scorer.
- The Manager or Coach of the home team, after being informed by the Official Scorer of the visiting team's starting line-up, must name the starting line-up of the home team. This information is then conveyed to the Coach of the visiting team by the Official Scorer.
- No changes in the starting line-up of either team, as given to the Official Scorer, or in the playing line-up on the ice, can be made unless approved by the Referee prior to the start of the game.
- Violation of this rule results in a bench minor penalty for the offending team. This must be brought to the Referee's attention prior to the second face-off in the game.

# Many Libraries One LLM

Nicholas  
Carr's *The  
Shallows*

What are some of the downsides of E-books and digital readers



## SOURCES

Based on the provided text, some of the downsides of E-books and digital readers mentioned are:

- The Net seizes our attention only to scatter it, making it harder to focus and read deeply (Chapter 7, Chunk 1)
- The dominant mode of reading is becoming skimming, rather than in-depth reading, which can lead to a lack of understanding and retention of information (Chapter 7, Chunk 7)
- Research shows that people who read linear text comprehend more, remember more, and learn more than those who read text peppered with links (Chapter 7, Chunk 3)
- The absorbed and personal mode of reading is discouraged when reading on digital devices, and readers' attention is directed towards the machinery of the hypertext and its functions rather than the content (Chapter 7, Chunk 3)
- Digital readers can be distracting, with features such as links, navigation menus, and notifications competing for the reader's attention (Chapter 6, Chunk 0)
- The linearity of printed books is shattered when reading on digital devices, along with the calm attentiveness it encourages in the reader (Chapter 6, Chunk 1)
- The way we read e-books will be very different from the way we read printed editions, which can lead to a loss of intimacy and connection with the text (Chapter 6, Chunk 1)
- Digital readers may lead to a decline in the art of writing, as authors and publishers adapt to readers' new habits and expectations (Chapter 6, Chunk 1)

These downsides highlight the potential challenges and changes that come with reading and writing in a digital age.




# Many Libraries One LLM

Dr.  
Lembke's  
Class  
Lectures



What is scheduling in an operating system like



 SOURCES

Scheduling in an operating system refers to the way the CPU time is allocated to different processes. It involves deciding which process should be executed next and for how long. This decision is made by the operating system's scheduler. The scheduling algorithms can be divided into four categories based on two quantities: preemption vs non-preemption and arrival time vs CPU need time. The main goals of scheduling are to keep the CPU busy, ensure fairness, and maximize throughput while minimizing turnaround time for CPU intensive processes.



# Why Rosie?

- Mixtral-Mistral8x-7B (The model we have the best results with) at 16 bit floating point precision runs well on four DGX GPUs (only half of the context window as well)
- This equates to 128 GB of VRAM. VRAM is basically how much data a gpu (fancy computer) can work with at one point
- A RTX-4090 (The best consumer grade GPU) has 24 GB of VRAM. It simply cannot load this size model
- A 7B parameter model at 16 bit floating point precision can be run on two T4s, still too much for someone with a RTX-4090
- RTX-4090 could quantize to 4-bits, however results are poop
- Why no quantize → VLLM does not have stable version of AWQ available
- Over 1000 hours of jobs run on chatter alone over the course of the last 6 months. Most of those hours coming in the last three months.
- We would have been running a p3.16xlarge instance on AWS, which costs 24.48\$ an hour, bringing us to a total cost of over \$24,000 to host this service on AWS
- Lets say you just want to run the library
- Library is loaded on two T4s, this involves loading several models on to it
  - Embedding model
  - Reranker
  - Library
- This could be run locally, but every model in use would have to be much smaller, the embedding model, avoid using the reranker model and run a quantized 7B model or lower
  - Quality would simply be lacking
  - Processing of concurrent requests
- One nice added unplanned bonus feature is that Various libraries can use the same LLM. If one LLM is stood up, various people could make different libraries on their own documents and use the central LLM.
- Why do we need the H100s though? We can use bigger and better embedding models. Load bigger LLMs

Mixtral-Mistral8x-7B: 23 on Chatbot Arena  
Mistral-7B-Instruct-v0.2: 37 on Chatbot Arena  
Llama-2-7b-chat: 51 on chatbot arena

## Estimates

- 12.24/hr With AWS Sagemaker Serverless Inference, 4 V100s
- High estimate for us: 1750W (½ max power draw of dgx1) when inferencing, 500W when idle
- For 1 day, need to be inferencing for 12 minutes to be cheaper than Sagemaker
- For 30 mins/day - Rosie: \$2.40, SM: \$6.12
- For 2 hours/day - Rosie: \$2.76, SM: \$24.48

## Amazon EC2 P3 instance product details

Instance Size	GPUs - Tesla V100	GPU Peer to Peer	GPU Memory (GB)	vCPUs	Memory (GB)	Network Bandwidth	EBS Bandwidth	On- Demand Price/hr*
p3.2xlarge	1	N/A	16	8	61	Up to 10 Gbps	1.5 Gbps	\$3.06
p3.8xlarge	4	NVLink	64	32	244	10 Gbps	7 Gbps	\$12.24
p3.16xlarge	8	NVLink	128	64	488	25 Gbps	14 Gbps	\$24.48
p3dn.24xlarge	8	NVLink	256	96	768	100 Gbps	19 Gbps	\$31.218

## Other Reasons

- Allows for very easy configuration and re-configuration
- Gives easy access to all logs, hardware config, and data inputs and outputs
- Little extra cost for data storage/moving

# Why Rosie?

- Big LLMs require big hardware
- Mixtral-Mistral8x-7B at ½ precision requires 128 GB of VRAM
  - 4090 has 24GB



# Why Rosie?

Our cost on Rosie:

- 3500W under full load
- 500W when idle
- Electricity cost: 19 cents/kWh

At 8 hours/day for 1 week: < \$50

AWS with "Sagemaker Serverless Inference":

- \$24.48/hour

At 8 hours/day for 1 week: **\$1370**

## Amazon EC2 P3 instance product details

Instance Size	GPUs - Tesla V100	GPU Peer to Peer	GPU Memory (GB)	vCPUs	Memory (GB)	Network Bandwidth	EBS Bandwidth	On- Demand Price/hr*
p3.2xlarge	1	N/A	16	8	61	Up to 10 Gbps	1.5 Gbps	\$3.06
p3.8xlarge	4	NVLink	64	32	244	10 Gbps	7 Gbps	\$12.24
p3.16xlarge	8	NVLink	128	64	488	25 Gbps	14 Gbps	\$24.48
p3dn.24xlarge	8	NVLink	256	96	768	100 Gbps	19 Gbps	\$31.218



# V100 Specs



## SYSTEM SPECIFICATIONS

GPUs	<b>8X Tesla V100</b>
Performance (Mixed Precision)	<b>1 petaFLOPS</b>
GPU Memory	<b>256 GB total system</b>
CPU	<b>Dual 20-Core Intel Xeon E5-2698 v4 2.2 GHz</b>
NVIDIA CUDA® Cores	<b>40,960</b>
NVIDIA Tensor Cores (on V100 based systems)	<b>5,120</b>
Power Requirements	<b>3,500 W</b>
System Memory	<b>512 GB 2,133 MHz DDR4 RDIMM</b>
Storage	<b>4X 1.92 TB SSD RAID 0</b>
Network	<b>Dual 10 GbE, 4 IB EDR</b>
Operating System	<b>Canonical Ubuntu, Red Hat Enterprise Linux</b>
System Weight	<b>134 lbs</b>
System Dimensions	<b>866 D x 444 W x 131 H (mm)</b>
Packing Dimensions	<b>1,180 D x 730 W x 284 H (mm)</b>
Operating Temperature Range	<b>5–35 °C</b>

# Other Advantages

- Allows for very easy configuration and re-configuration (yaml)
- Full control over logs, hardware, job time, LLM, embedding model,
- Ingress-Egress fees are close to nonexistent

- It's all on-site



```
sbatch_job_time: "0-12:0"

llm_config:
  folder: "new_backend/LLM"
  engine: "vllm"
  partition: "dgx"
  requirements_file: "llm_requirements.txt"
  gpus: 4
  cpus_per_gpu: 16
  max_model_len: 24000
  hf_llm_name: "mistralai/Mixtral-8x7B-Instruct-v0.1"
  dtype: "float16"
  out_file: "llm.out"

db_config:
  folder: "new_backend/db_guardrails"
  partition: "teaching"
  txt_dir: "/data/sdp/senior_design_llm/policies_v3/"
  gpus: 2
  cpus_per_gpu: 4
  reload_db: true
  reranker_name: "BAAI/bge-reranker-base"
  requirements_file: "backend_requirements.txt"
  embedding_model_name: "BAAI/bge-large-en-v1.5"
  prompt_file: "config/system_prompts.yml"
  metadata_file: "/data/sdp/senior_design_llm/policies_v3/metadata_generated.json"
  out_file: "db.out"
  summary_chunk_size: 512
  summary_chunk_overlap: 64
  document_chunk_size: 1024
  document_chunk_overlap: 128

guardrails_config:
  folder: "db_guardrails"
  reload_guardrails: true
  cached_guardrails: "cache/guardrails.yml"
  guardrails_config_file: "config/guardrails.yml"
  embedding_model_name: "BAAI/bge-large-en-v1.5"
```



- Four dgx GPUs running at float16 precision for our LLM.
  - Mixtral 8x7b
- In our pipeline to retrieve information there are 3 models at play
  - Embedding model from BAAI
  - Reranker model from BAAI
  - and an LLM
- Each taking up multiple Gb of VRAM.
- Calculating Embeddings, rerankingrunning LLM inference and reranking
- Our pipeline takes around 4-8 seconds to load for most documents. Asking questions takes a very long time.

# Self healing Huberman Podcast

- 1650 documents
- 3 questions each
- 4950 questions
- 5 seconds each → This is one library visit
- 24750 seconds
- 412.5 minutes
- 6.875 hours per round



## To use this?

- Simply run git clone on our repo
- Copy and edit configuration files to your liking
- Run two python scripts and wait
- There will be a detailed guide soon... (With pictures and everything)