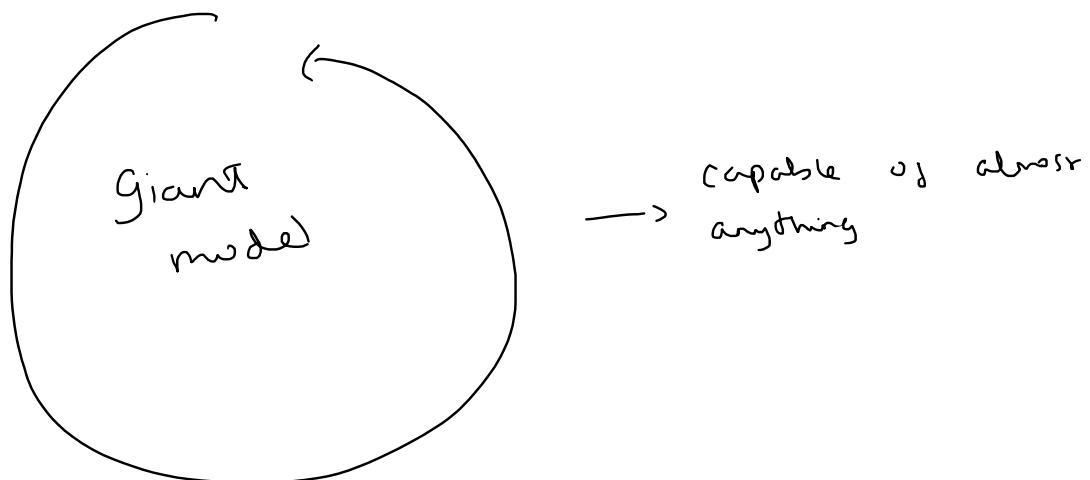


Gen AI agents have, until now, mostly been built using brute force approaches. That is, massive amounts of data are passed in and the models refine themselves.

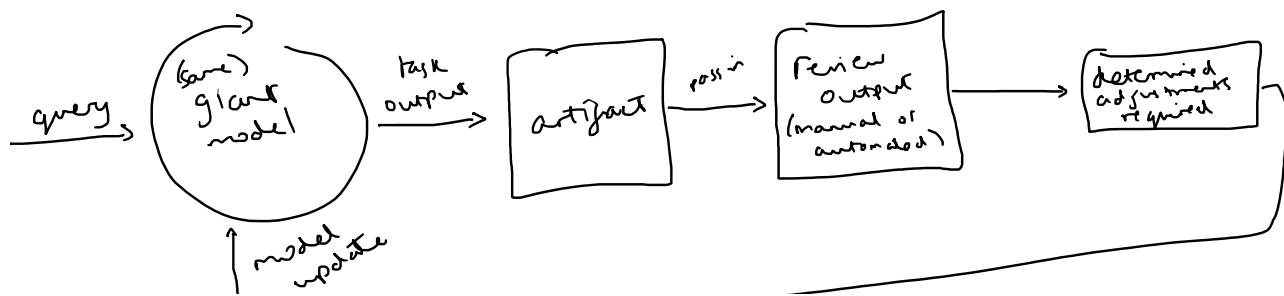
This should work to begin with, or it has, but will plateau based on how well those agents can refine themselves to ensure adequate capability and understanding of all required science and technologies.

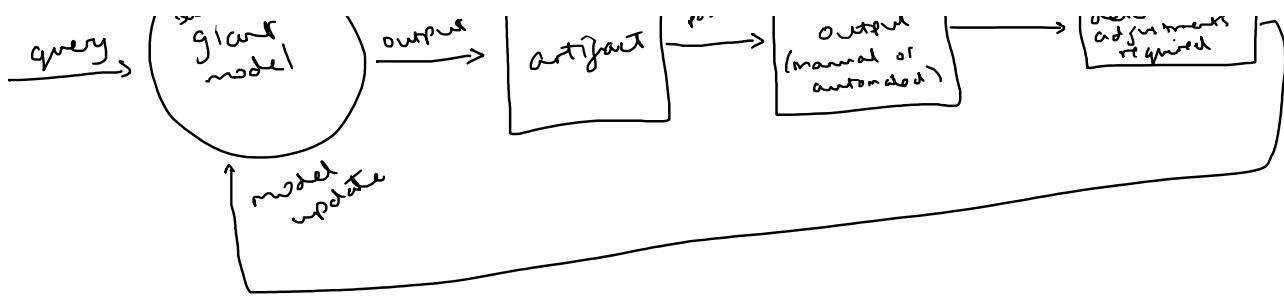
Because the systems were created by brute force they will be starting as massive data units. The next steps may present themselves as reduction tasks but they actually become construction tasks.

Starting point:

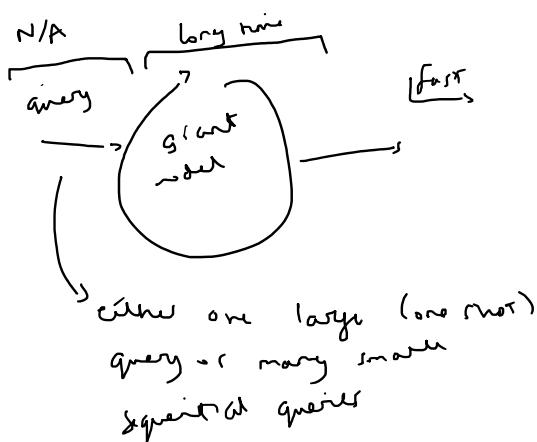


Red herring technique:



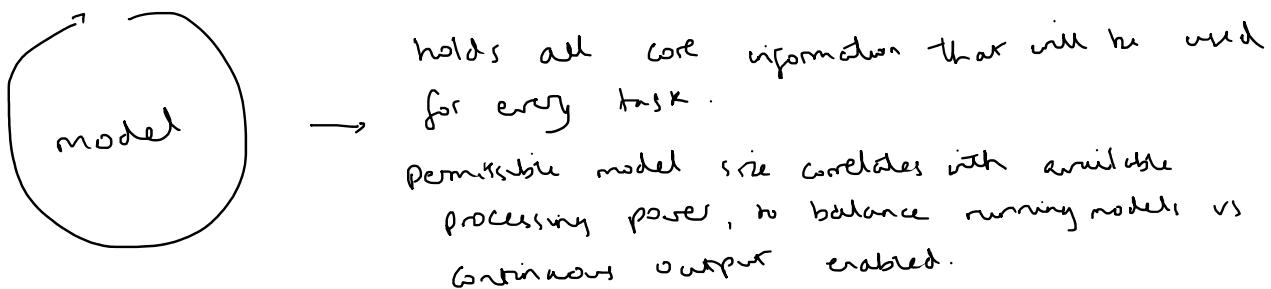


This is a red herring (false lead) because the overall model is simple and the noted time complexity would soon promote how the wall time will be unnecessarily large owing to the single model's size.

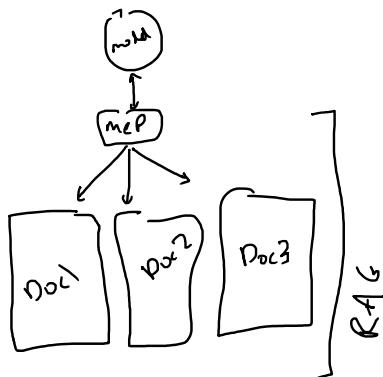


An inadequate artifact requires complete re-creation. Training a different model from the ground up would be results very slowly, such that it will still appear more fruitful to refine the large model but that equally ignores the value presented by other adjacent technologies.

Reference: "we're to handle "what"



If processing power is variable then multiple model sizes may be used, assuming that same scaling rule is available for other components too.



Holds all the per-use-cache information that is only stored depending on the task at hand.

→ Document sizes depend on the available storage space as well as available cache (or cache equivalent) for the processor.

Beyond size, documents should be chunked based on their purpose and "skill tier". e.g. CS 101 and CS 421 shouldn't be in the same document. CS 101 and CS 102 could be better than CS 101. There is a chance that CS 101 is discardable in most situations. (There is a lot more to discuss on this.)

Documents in modern RTG systems are handled as an almost deterministic search problem, but that is not correct.

Retrieval is also viewed as a lower problem which is also not correct.

### Side note: determinism

The main limitation in digital GenAI is the lack of determinism. Because every step from query processing to RTG retrieval is non-deterministic the overall system behavior makes destabilization after even the smallest model update.

An unwritten goal is (or should be) to make as much functionality as deterministic as possible.

This will make the system appear more fractured owing to the numerous additional models and steps in system diagrams.

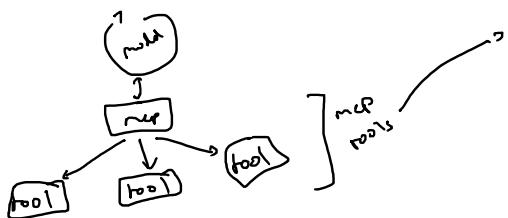
It will however create a system that is more maintainable, scalable, transportable, and - most importantly - reliable.

The reason some countries have allowed themselves to counts with enormous populations is because it is easier to train a lot of instances to correctly perform a smaller task than it is to train a few people to reliably perform entire sets of tasks.

A larger population focuses on individual task optimization, whereas a smaller population focuses on maximizing each individual's ability.

These same principles translate directly to Gen AI.

Lemma: Reference



MCP tools are the rules for most deterministic behaviors.

In fact, MCP tools should be the rules for all semantically defined repeated behaviors, but most importantly for all deterministic behaviors.

MCPs are modular by nature and should remain as such - any black box hidden logic should be kept to an absolute minimum.

grouping functionality within single MCP tools is convenient but, as per the discussion on population sizes, we should allow a resource intense end state and favour logical separation of MCP units at every logical boundary.

task composition is the end state and MCPs should be prioritized as such.

The most logical boundary for MCPs is based on which data is required. This enables optimal data storage and use in any final state system.

Sidenote: MCPs are not long lived processes

(at least, not always)

While MCPs are generally presented as permanent system units (always on or always off), they will actually be launched and kept running based on real time demand.

That is harder to see right now because processor availability far outweighs total demand.

As modular deterministic behavior becomes more available, the barrier to entry gets lower and lower.

As well, scale and complexity currently bounce back and forth due to the naturally coupled nature. In the future scale will increase in parallel with complexity improvements.

That is: in the future, the teams adding scale to applications will have no knowledge of what the deterministic tools are doing or how well they perform. Processor demand will, at some point, increase almost exponentially and real time efficiency will be more important than ever before - it will happen seemingly overnight.



### Sidensle: multiple models

While a brute force approach at present will begin with a single (or few) large models, there will inevitably be multiple models in use too.

Right now models are mostly compared based on cost per unit with little regard towards other equally important contributing factors. Things like query time resolution are abstracted to "best or your pricing tier" with little discussion on why variance is introduced. This assumes things that would dictate a need to split one model into two.

### Resume: Reference

The narrow assumption is that models should be


 → Split based on the roles they will do.  
that is not true!  
 because processing becomes the main constraint,  
 models should be split based on data boundaries - not  
 the data they will use, but the data they need  
 to understand.  
 in that sense, you might group one model for  
 chemistry, one model for biology, and one model  
 containing physics and maths which have  
 significant overlap.  
 in a RAG system, any details about specific  
 programming languages should be stored as  
 optional data. the only core understanding  
 required by a model should be how  
to program in a way that translates to any  
 language. this aligns with how the most prominent  
 programmers operate too, since core principles will  
 stand the test of time while floating language choice  
 may not.

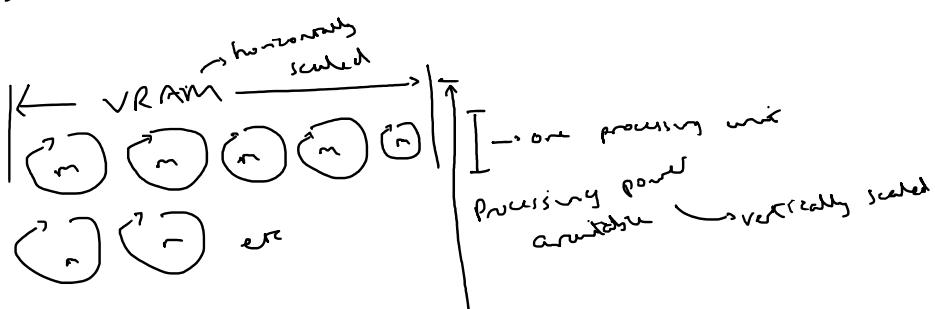
## Reference: Scaling

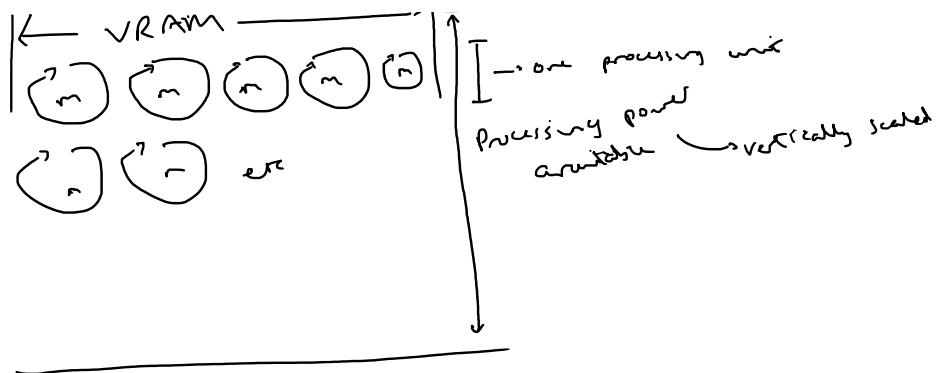
Scaling model use is optimally achieved when the models correctly fit the hardware, either exactly or in some multiple amount.

The scaling is achieved in two dimensions:

1. Processing power; and
2. VRAM required

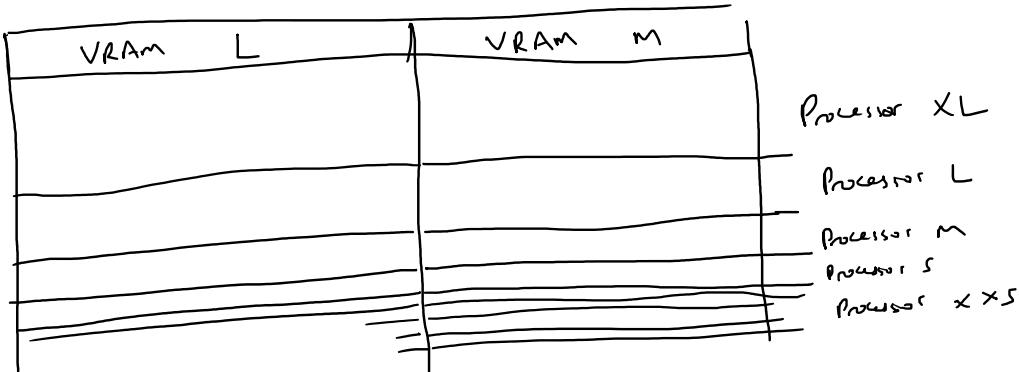
The two may be related but it is a case of how much time is spent "doing"  
 vs how time is spent gathering or waiting.





In practice, the running models will not be so easy to predict ahead of time. That is when other hardware capabilities take priority during algorithm choices. As well, VRAM is more likely to be constant than processing power is more likely to differ per processing unit available.

The actual grid will look more like:



In this scenario, high single-units of processing power are rarest. So, models will be run using a greedy algorithm to "run larger models first using the least populated required processor". VRAM will only be promoted if a request to run a model for a real time request.

### Reference: Direct Storage APIs

Direct storage was advertised as something to enable unbelievably rapid game load times. It quickly fell silent given:

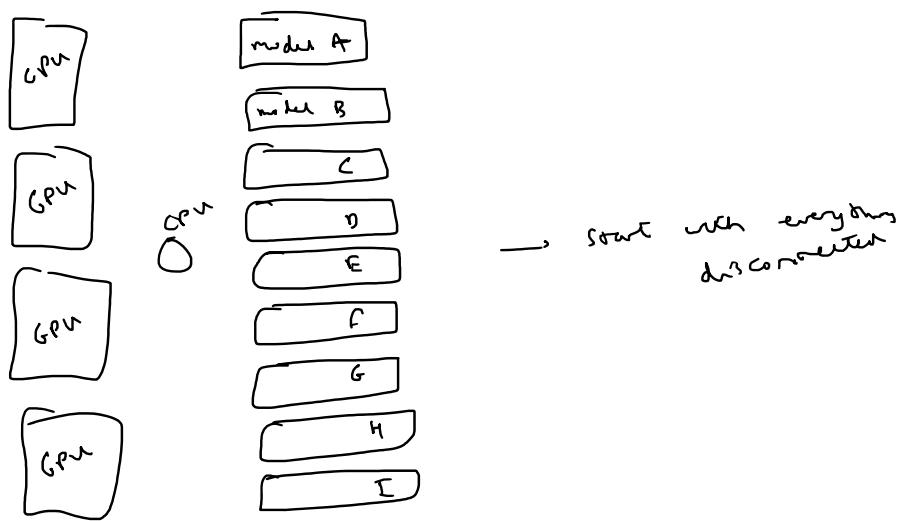
1. Implementation cost outweighed immediate adoption benefit; and
- This organization would be abstracted out to

- Implementation benefits; and
- This optimization would be abstracted out to game engines in the long run anyway.

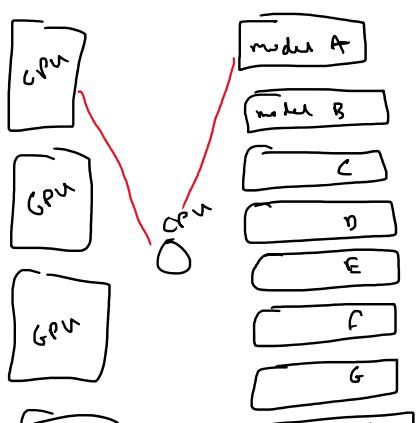
It has not (yet?) resulted with reference to the unique power it should afford GenAI.

The benefit in GenAI is essentially hot-swapping models. Even in systems where the processing units are isolated, efficient storage switching will enable models to be whatever processing power is available without interrupting any running process.

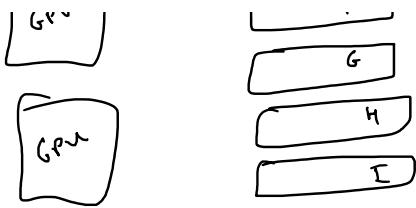
Even in the case of a single total system, it will allow pseudo infinite storage limited only by available storage hardware, regardless of any drive connection limits.



→ Request to run "query" on model A



1. Connect model A storage unit
2. Direct load model A to GPU
3. Once model A is on GPU, unmap data if not set to auto evict  
say to disconnect model A storage unit again.



It's not realistic that all storage will begin unplugged. The most commonly used models will likely be converted in anticipation with options to swap out as needed in real time.

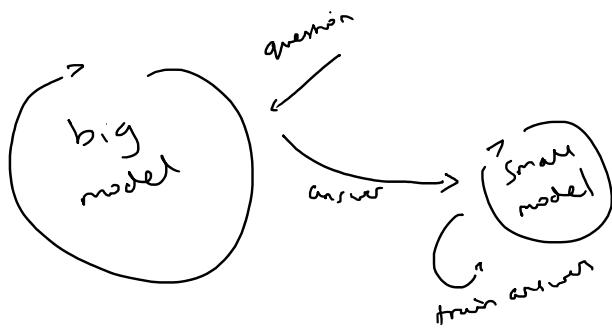
### Reference: model growth

Since reducing a large model to the core components is not true efficient, it is preferable to refine model understanding using smaller models and then transfer complete learnings to larger models.

This goes against the grain of model distillation because we are adding a priority processing time (both wall time and time spent on chip).

This will be the fastest way to achieve maximum performance models using the minimum processing power.

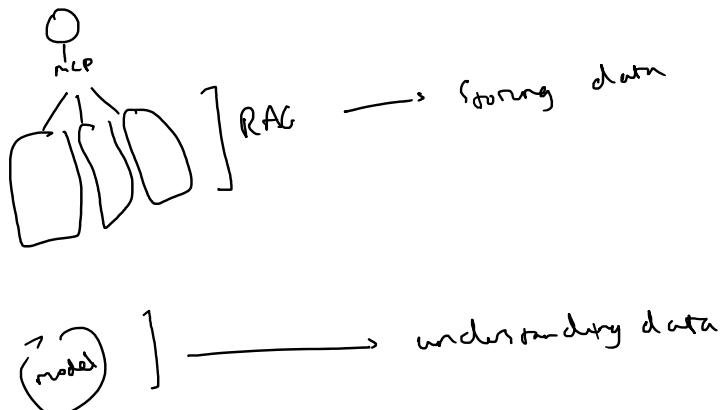
### Model distillation:



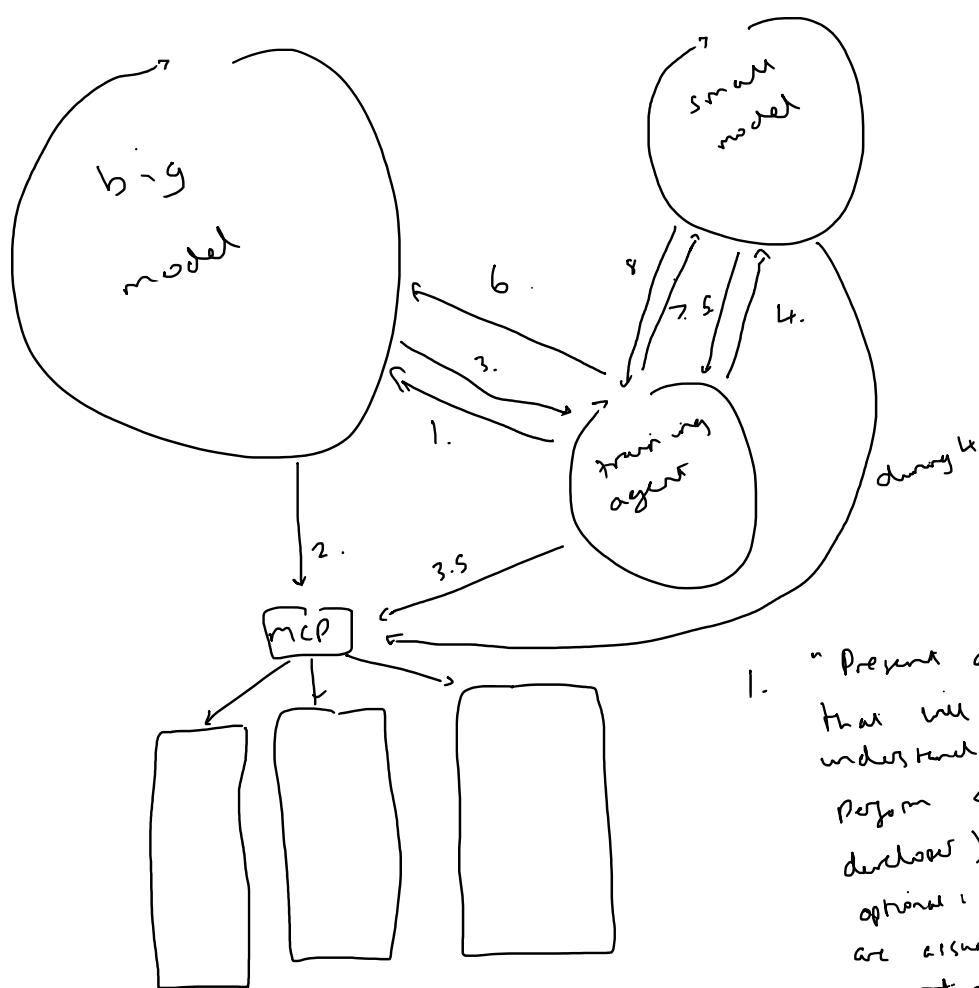
and the small model grows over time

Model distillation completely ignores the presence and utility of adjacent technologies like RAG and mLP. It assumes an absolute need to grow smaller models since more efficient data storage may suffice.

To reiterate from earlier reference on component purpose:



So optimal model growth should actually handle like:



1. "Present a complete question-answer turn that will test for a thorough ability to understand all documents required to perform <role>" (eg role = L2 python developer).  
optional: also specify which documents you are asking will be accessed for such operations

2. Query RAG for data on specified topic ; create test

2. Query RAG for data on specified topic ;  
create test
3. Pass test to training agent  
optionally include: document list
- 3.5. (Optional) Store cached requirements (and optional document list) in case of stability in future training tasks (will review example later).

- loop these two steps
- if and only if
- employing option b.c.
- If option b.c. is employed  
the test should be written on the basis of existing and known additional documents found during the substitutes found during the substitution search.
4. Administer test , allowing access to RAG system for documents

option a: no document limit specified

option b: limited to only any specified document list

option c: either a or b with an option to later review additional document in case of wrong answer.

5. Provide test results

6. Based on final test results, query for training data to fill identified comprehension gaps.

7. ] Practically equivalent to 4 and 5
8. ]

optionally b.c.

The overall goal is to minimize uncertain calls do the largest (or most expensive) model.

(Creating and caching a question-answer test enables reuse when training other models, especially if the test is appropriately partitioned per topic.)

e.g. If CS 101, 102, and 103 are all tested and cached and then we need to train a biology model with preq. CS 102 we can just use that cached test instead of asking again.

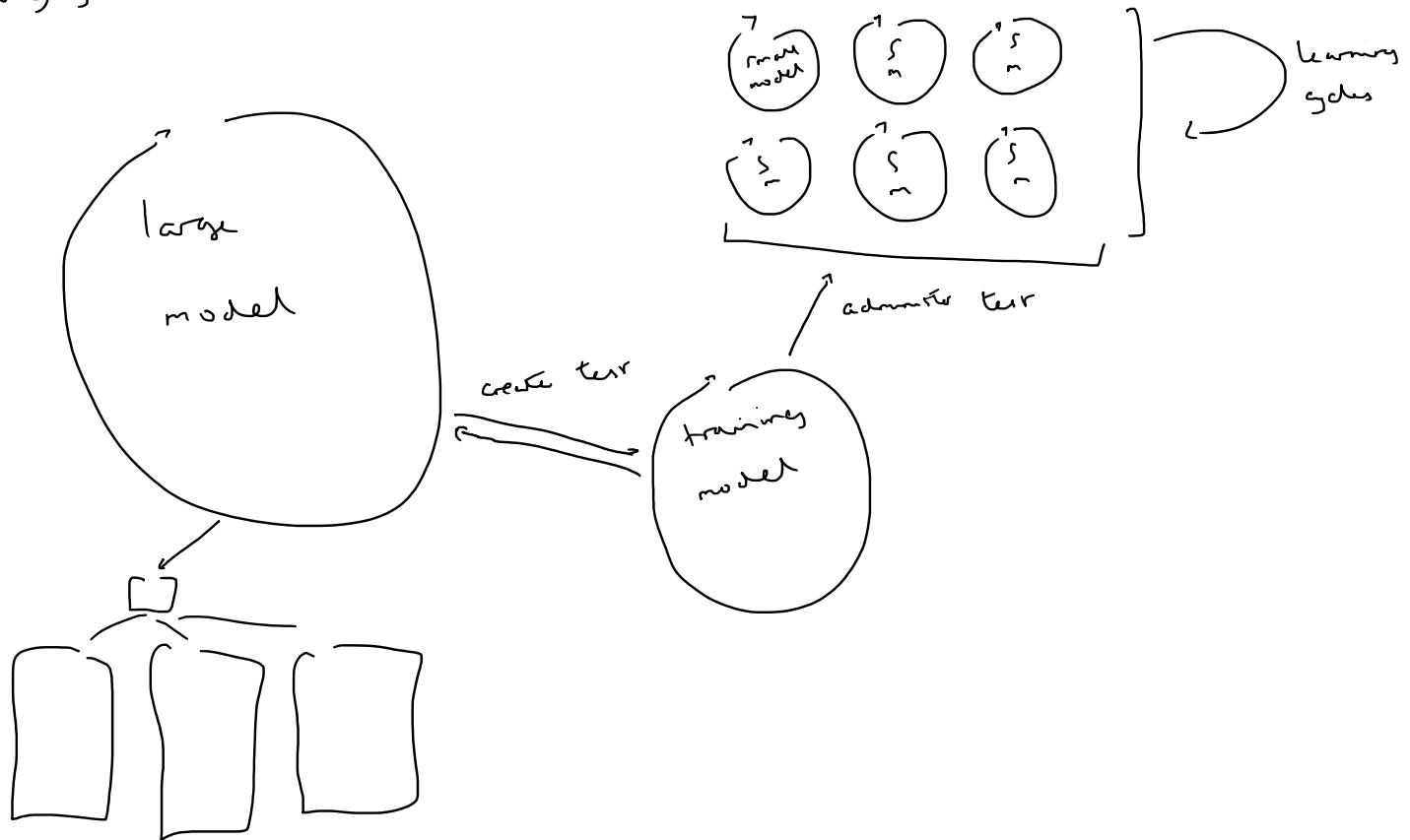
We can still optionally check that any document list returned is

still complete and unchanged.

An downside is that the training agent creates the questions and only gets complete answers from the large model. This assumes understanding the data is not required to further query it, only to provide concrete answers.

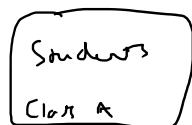
As well, in case of employing LCF where the student model can investigate for further documents, it presents an opportunity to prioritise work time.

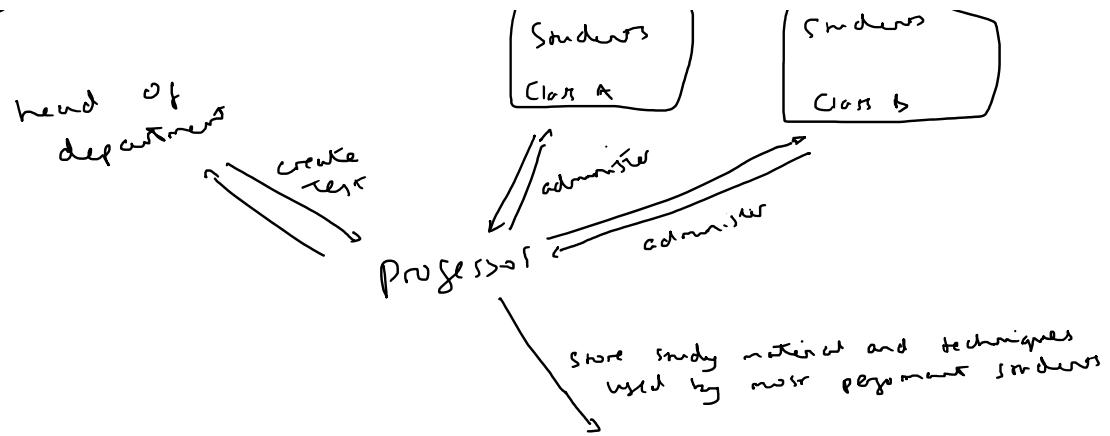
If ample processing power is available for the smaller models, multiple models can run in parallel each employing functionally different "studying" techniques



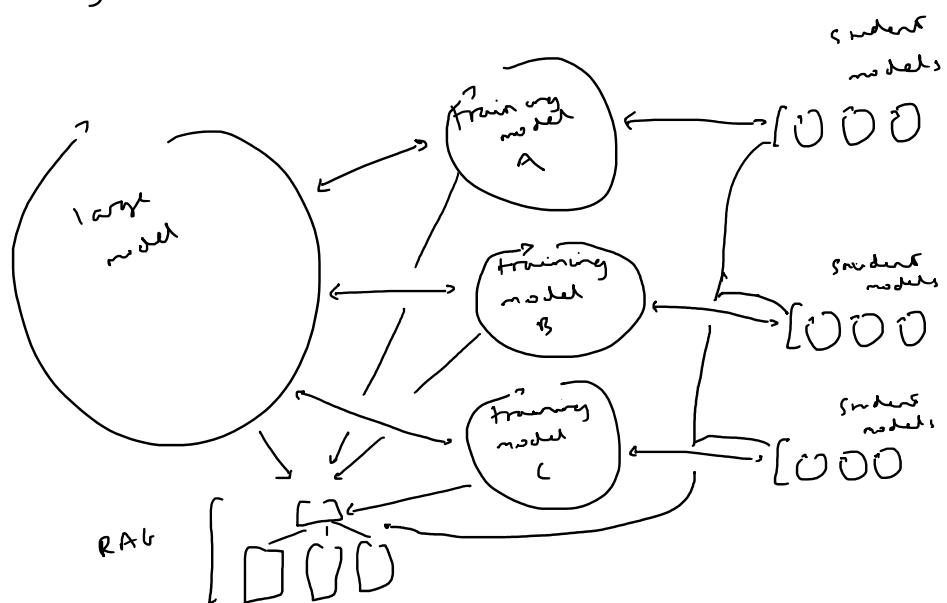
This is following the same rough model used in high schools and colleges

, ... of

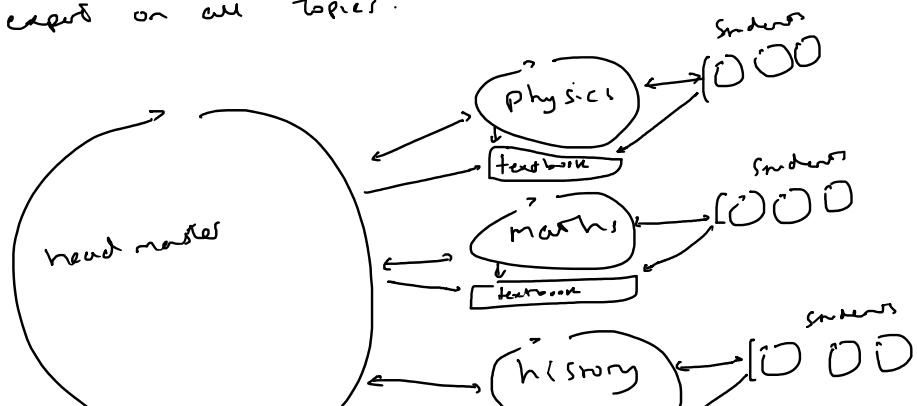


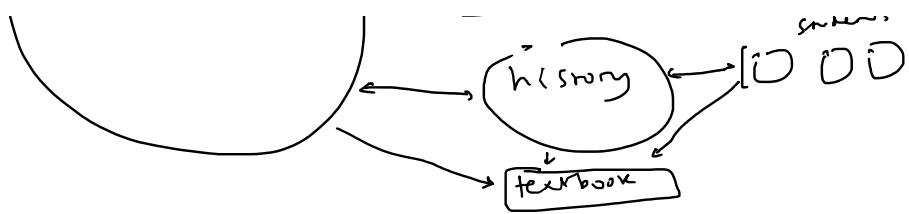


Printing the time spent by the largest model also serves to enable teaching the most variety in the shorter time. This is done by scaling the training model and student models, which will have hardware most easily available anyway.



This follows a slightly older paradigm where there would be subject teachers each reporting to the headmaster who is an expert on all topics.





Reference: Cross learning

Reference: Data intense application