# Generative LLMs: How they work and where they are headed

Eduard Hovy

University of Melbourne and CMU

**Look what an LLM can do!**

*Why can it do that?*

— I have no idea / that's future work / I've never thought about it

**Look what an LLM cannot do!**

*Why not?*

— I have no idea / that's future work / I've never thought about it

**I don't care about LLMs — here's what I did**

*Why are you doing this? Can't an LLM do it already?*

— I have no idea / that's future work / I've never thought about it
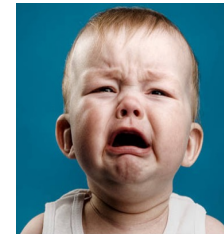
# NLP fundamentally changed direction last year

## Old NLP, in 2035

- **Research**: Many 'smaller' new areas like social roles, extremism, style, mis/disinfo, author/user personality, etc. LLM mining, but no sign of GenAI

- **Engineering**: Tuning LLM algorithms and datasets to larger applications (QA, IE, MT, IR, text mining, dialogue, etc.)

- **Evaluations**: Increasingly focused (both black-box and glass-box), some new automated metrics

- **Use**: Large corporate frameworks that integrate NLP and multimedia tasks and sell service on the cloud

## New NLP, in 2035

- **Research: 'NLP Lite'** with prompt programming using GenAI as the tool. Many more people do it, also non-researchers

- **Use**: Tailored GenAI services on the cloud

- **'Real' NLP** research — *but on what?*



**Are you sad or happy?**

# Three major directions for New NLP

**Research on LLMs**

**Research using LLMs**

**NLP engineering: Make LLMs <u>usable</u>**
- Build smaller and cheaper LLMs
- Systematize prompt engineering
- Integrate language, images, and other media and functions

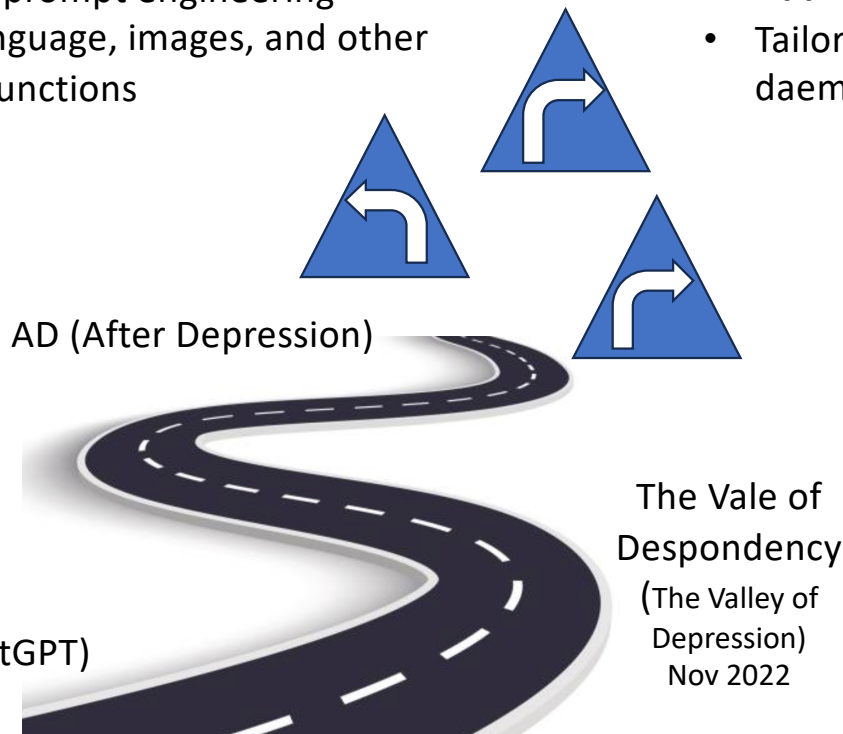**NLP applications: Make LLMs <u>useful</u>**
- Tune LLMs to domains and companies for enterprise processing
- Add functionality and agency in the world
- Tailor LLMs to people to be their personal daemons/amanuenses in everyday life

**NLP research: Make LLMs <u>understandable</u>**
(or at least, be solid engineering)
- Fix the problems with LLMs
- Get explanations how LLMs do what they do
- Formalize them well enough for autonomy, assurance, and ethics

**Some of both**

AD (After Depression)

BC (Before ChatGPT)

The Vale of Despondency (The Valley of Depression) Nov 2022

# 1. Usable

NLP engineering

# Making LLMs **usable**

## 1. LLM construction
- Build smaller and cheaper LLMs
- Enable LLM evolution (content addition)

## 2. Multimodality
- Merge text, voice, video, other behaviors
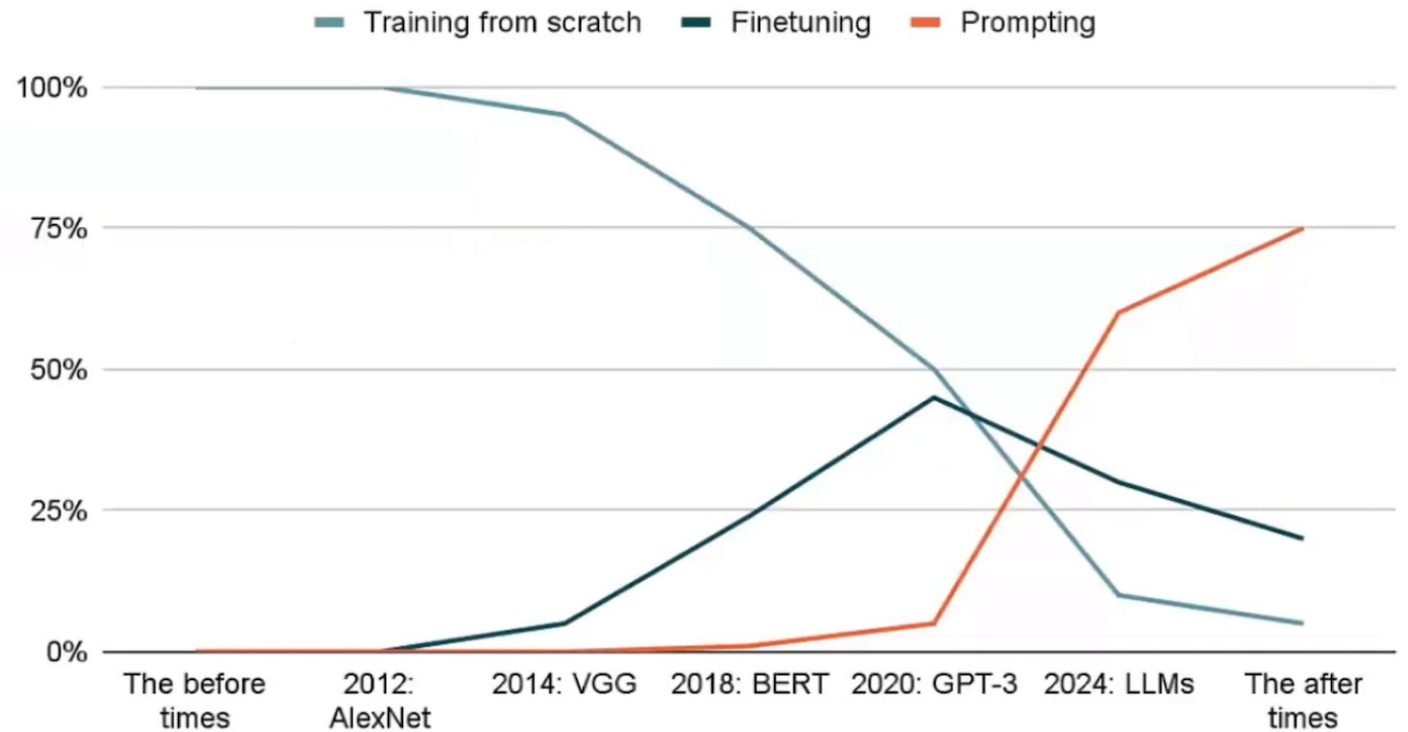- Enable true multi-modal dialogue

## 3. Prompts
- Systematize 'prompt space' and teach prompt engineering
- Create prompt evaluation metrics: how 'good' is a prompt? And its answer?

- Everyone wants an LLM but no-one trusts 'Open'AI or the other giants
- Smaller/weaker LLMs are acceptable for limited domains and tasks
- Owners want to fine-tune content: RAG etc. What tradeoff between model tuning/training and RAG?

- People know when to switch media in their stories — "which info do I need to draw, not say?"
- Encode embeddings from various media together, then process further
- Other modalities: movement, emotion…

# Prompts vs training

**Prompts are becoming more important than training**

Machine training vs human prompting

# Research topic 1: Prompt creation

- Systematize 'prompt space'

- Prompt tuning (Li and Liang, 2021; Lester et al., 2021)
    - Can you map specific task training info to the pre-training data?
    - Prompt tuning adds 'soft prompts' (embeddings) into prompts
    - NN must decode training label words in the same way as the pre-training objective

- Prompt writing
    - Multi-layer prompts with internal prompt structure: OpenAI's functions language inside prompts for hybrid English+programming: "take it step by step" / "if the answer is X then do A else do B"
    - Use prompt variations + adversarial models to discover prompting strategies
    - Develop teaching of prompt engineering

# Prompt engin**X**eering —> programming

**Prompt as a programming language**

- Quasi-English

- Non-deterministic

- Full of special-purpose keywords and tricks

**Execution engine**

- Deep neural net, not a normal computer

- Also non-deterministic and

- Constantly updated/changing

**Situation today**

- Still seen as "only" prompt "engineering"

- Already many online training courses and prompt-writing competitions

- This is what most people will use on graduating, *far* more than programmers

- But no serious attention from CS departments

# AUTOMAT for prompt writing

**Prompt component**                                          **Example**

- **A**ct as a —— : Bot persona and expertise               …journalist, expert

- **U**ser/reader persona and needs                         …for a 10-year old

- **T**argeted action to do: specific task                  …summarize, explain

- **O**utput description and format                         …table, bold headings

- **M**ode and style (formal, minimalist, etc.)             …legal, poem

- **A**typical cases / anomalies                            …for missing items, etc.

- **T**opic whitelist / exceptions: what to avoid           …only 2024, not 2023

# Research topic 2: Prompt evaluation

- What's the metric? What to evaluate?

- Metric desiderata:
  - **Shorter** is better — word count
  - **Accurate** is better (no hallucinations) — Precision on facts
  - **Complete** is better (no gaps) — Recall on facts
  - **Coherent** is better (proper answer structure) — discourse coherence metrics
  - **Understandable** is better (explanation included?) — dialogue task metrics

- Need a dataset of Prompt+Response pairs
- Will lead to automated prompt learning by optimizing against the metric
  - How do we determine theoretical maximum?
  - How can we extend prompt functionality?

# 2. Useful

NLP applications

# Making LLMs **useful**

1. Must contain the right info
   - Correct, complete, trustworthy
   - Knowledge must be updatable

   - How do you confirm that the content is good?
   - Time and change. How can you delete or change factoids? (May become a legal requirement)
   - RAG stores and updating methods

2. Tuned to domains and companies
   - Starter platforms and extensions
   - Multimedia technical knowledge ingestion

   - Build starter platforms (LLAMA-3 is popular)
   - Extend with *plug-ins* and the GPT Store
   - How can you read domain text/images/notations on top of that?
   - Specialist knowledge formats and reasoning

3. Tailored to people's lives
   - Privacy and IP protection
   - Personal style

   - Ingest a user's Facebook pages, photos, and personal docs on top of a generic User model
   - How do you organize info by timeline?
   - How can you implement (differential) privacy?
   - What's the economic model? How much will people pay for their own assistant in the cloud?

# 3. Understandable

Long-term NLP research

# Making LLMs **understandable**

## 1. Fix the problems
- Hallucination
- Truth recording and maintenance
- Timestamping and updating info

- Verify output: How to check for hallucination against the web?
- How to represent facts/truth?
- Long-doc output coherence: 'template' learning and editing

## 2. Make explanations (for users)
- Recognize and produce generalizations / rules in the LLM
- Record / find instances (training examples) for any situation

- Understand the iterative deepening nature of explanation
- How to tailor explanations to the user's knowledge?
- Make *Influence Function* approach practical and explainable

## 3. Understand LLMs (for researchers)
- Embeddings and microfeature processing
- What can they not do?—content and inference

- Probing methodology
- Guide the formation of kinds of generalized knowledge
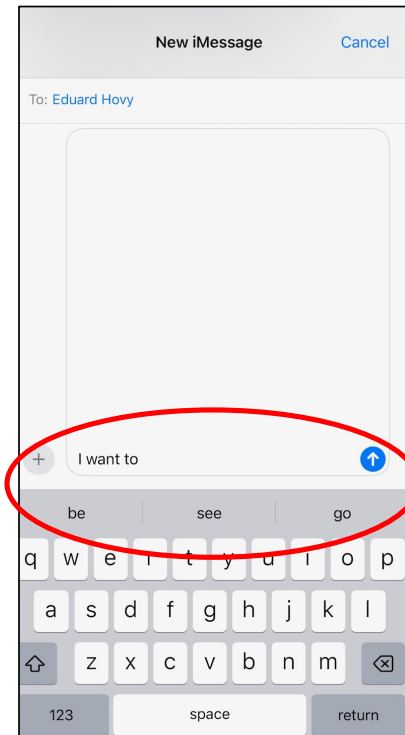- 'Microfeatures' and *reasoning circuits*

# Two deep challenges

- Understanding how reasoning / inference happens
- Generation: why is it so good?

# Why do LLMs seem able to reason?

# Do LLMs reason?

- A 'traditional' LLM like BERT is a **knowledge base** and an **info transformer**
  - It knows simple facts from statements
  - It can put together similar statements to find simple inferences

- A Generative LLM like ChatGPT adds a **chat loop** on top. Given a prompt it can continue concatenating to extend its output forever
  - It can put together sentence fragments that are true, false, or completely unjustified

- From the outside both look like reasoning

- When it makes mistakes, how can you understand what kind of error it is? You probably need to fix different things

# Language Models

Your cellphone can predict what you are about to type:

*I want to [go, eat, …]*
*I go to [the, a, her, …]*

How does it know this?

It needs a list of possible continuations:

"want    to    ?"

| w1 | w2 | w3 | score |
|------|------|------|----------|
| want | to | go | 0.000391 |
| want | to | eat | 0.000015 |
| … | | | |
| go | to | the | 0.000491 |

# Large Language Models



L="station" —> train-related
L="the" —> noun
L="to the" —> location
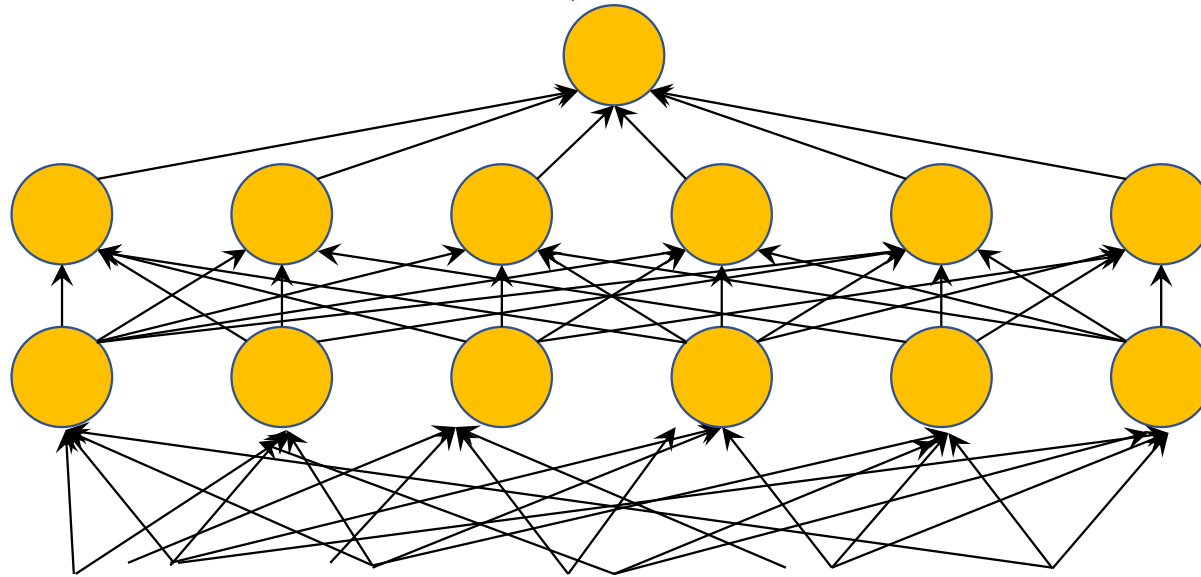L="on" —> location

"platform"

R="train" —> train-related
R="get on" —> near vehicle
R="9¾" —> "Platform"
R="wait for" —> location

Abstract
'microfeatures'

Weight scores on
each link adjusted
during learning:
backpropagation

"at the station he went upstairs to the _____ and waited for the train"
"At King's Cross Station on _____ 9¾ you can get the Hogwarts Express"

Hundreds of millions
of sentences

# LLM as a 'lookup table'

A big magic multi-D table of English

You give it [sets of] sub-features in given positions …

… it gives you the sub-features in the other positions

*When was the Panama Canal built?*

give me a time word

near to "Panama Canal"

near to any word(s) meaning "build"

L="when" —> time expression
L="Panama Canal" —> relating to
the Panama Canal
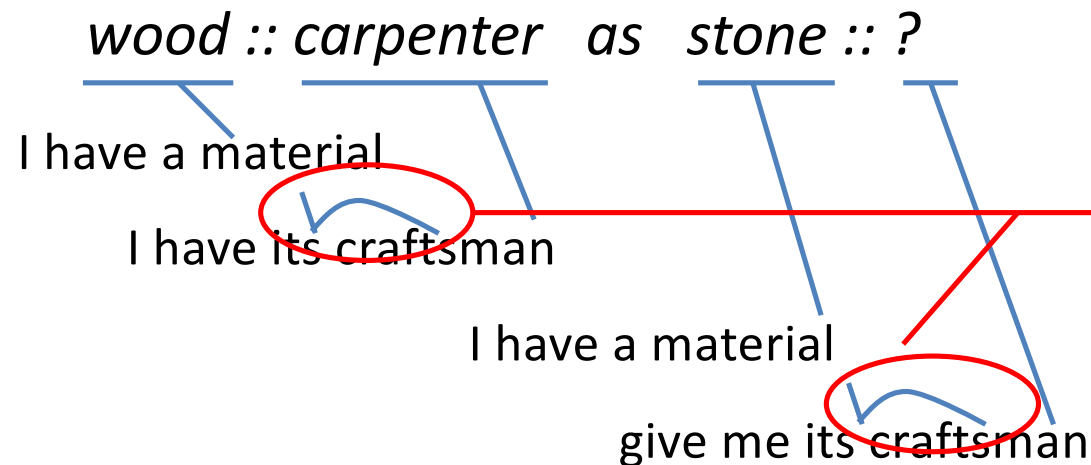L="built" —> construction
L="?" —> "the answer is <ANS>"

The answer is 1904–1914

# LLM as a 'lookup table'

A big magic multi-D table of English

You give it [sets of] sub-features in given positions …

… it gives you the sub-features in the other positions

*wood :: carpenter   as   stone :: ?*

I have a material

I have its craftsman

I have a material

give me its craftsman

This microfeature connectivity graph is like a programming language

# Skill neurons

- After pre-training and before tuning, some neurons are reliable predictors for task performance. These neurons are:
  - stable across trials
  - required: if they are perturbed performance drops significantly
  - task-specific: similar tasks exhibit similar skill neuron sets, and skill neurons of same-type tasks are more important for handling a task than those of different-type tasks
  - not activated shallowly from task keywords: their predictions are not significantly influenced by the label words used in prompt tuning
- Can use them to prune NN — keeping just the top skill neurons active, can reduce 1/3 of pre-trained Transformer parameters and get about 1.4x inference speedup

# Some approaches

1. Knowledge: Can we edit LMMs and change what they know?
   Wang and Ji 24, Wu et al. 24, etc. from UIUC

2. Inference: Can we delete/ablate specific reasoning paths/circuits?
   Overview by Miller et al. (COLM 24)

3. Reasoning chains: Can we identify paths that have specific effects?

# 3. Identifying individual reasoning paths

Questions:

- How can we find a reasoning 'pathway' for a *specific* task computation?
- How can we display it?
- Is the pathway we find truly doing what the NN does overall?

Some relevant work

- Expert nodes ("skill neurons") — Wang et al., EMNLP 2020
- Information Flow Routes (IFR) — Ferrando & Voita, arXiv 2024
- Anthropic's 'circuits' and the Golden Gate Bridge example
- Identify subparts of the NN as a functional 'circuit'/pathway that computes a specified portion of the output

# Our method: contrastive semantics

There's too much noise!

Need to identify a <u>minimal</u> inferential difference between two prompts

- Example: Winograd sentences  (Winograd 1972)
  - Two sentences that differ in only one or two words and that contain an ambiguity that is resolved in opposite ways, and requires world knowledge and reasoning to be resolved

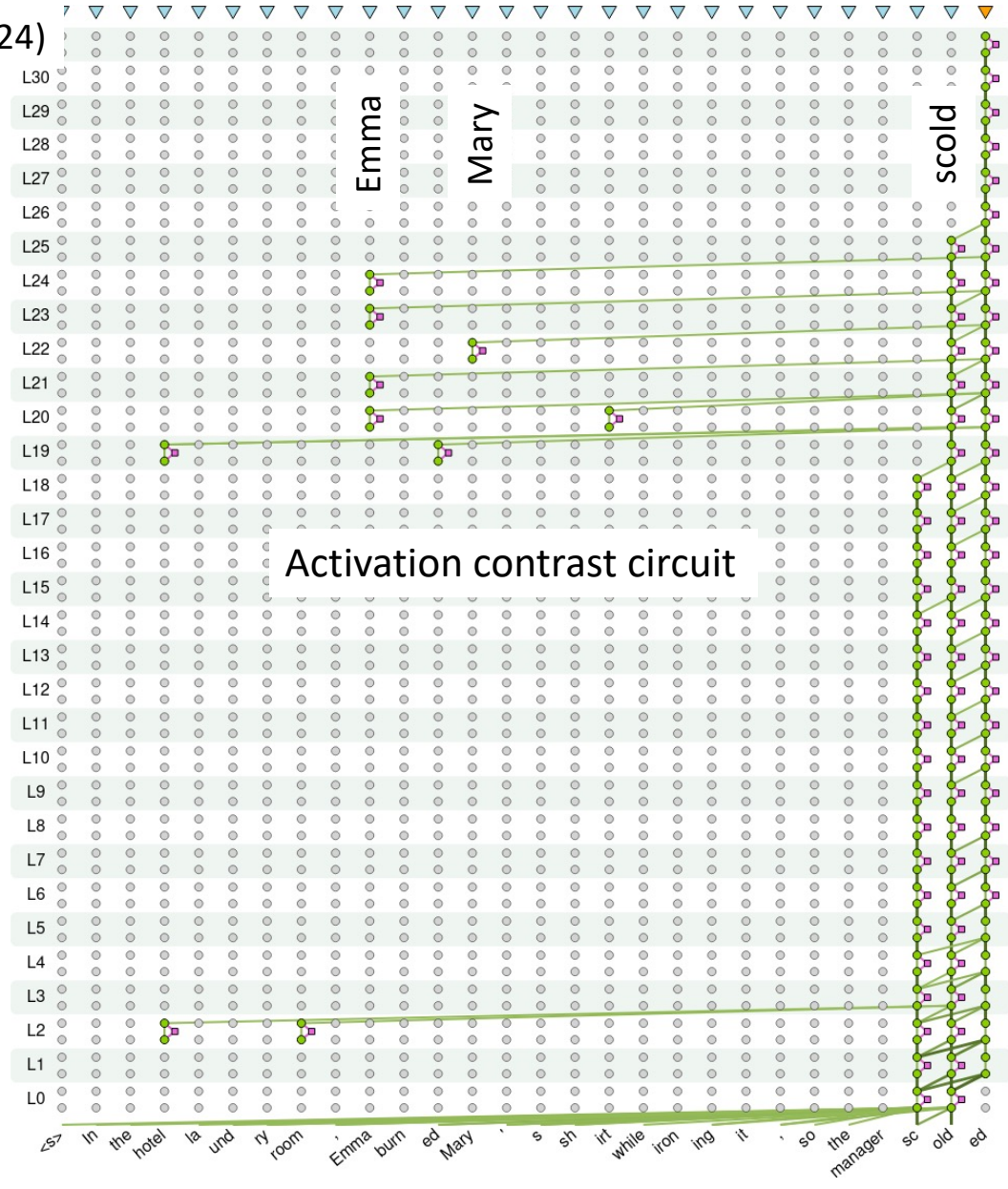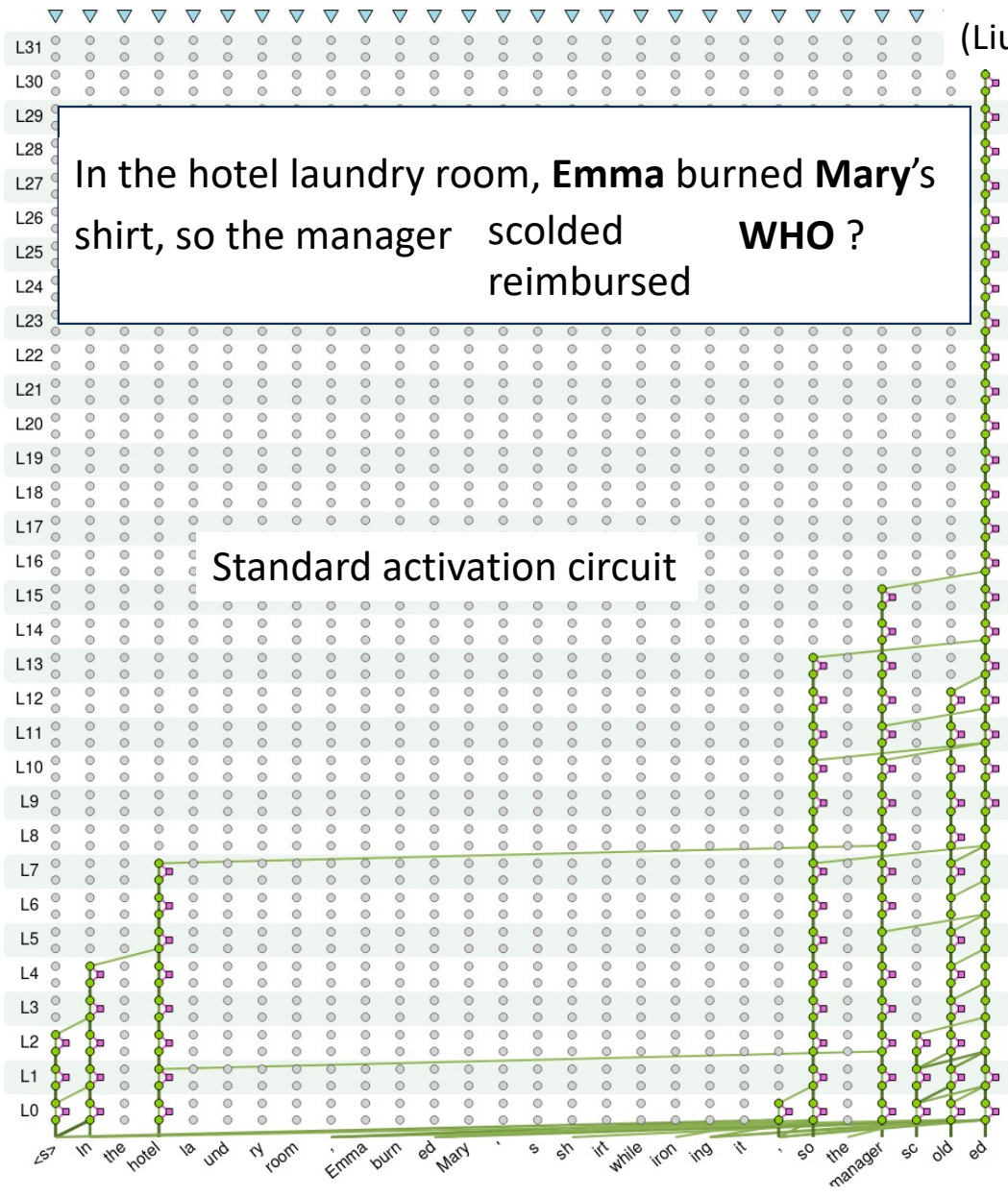> The city council refused to give the demonstrators a permit because they [feared/advocated] violence

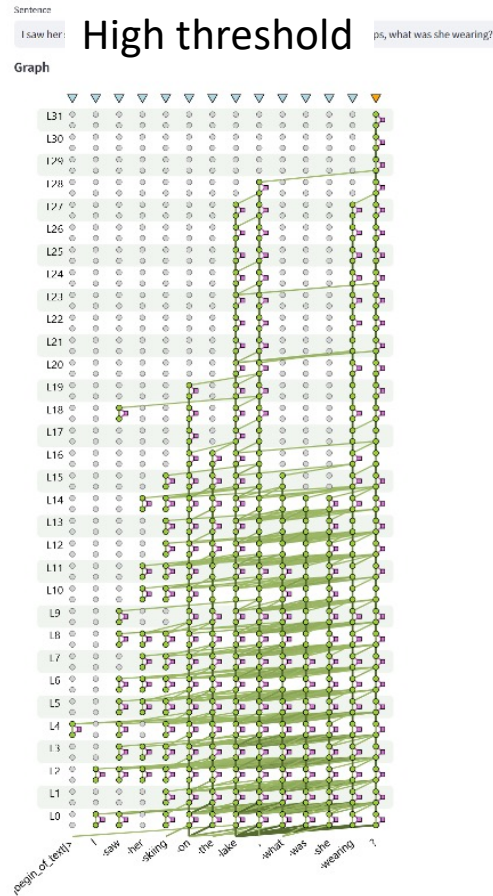Who is "they"?

"feared" —> city council
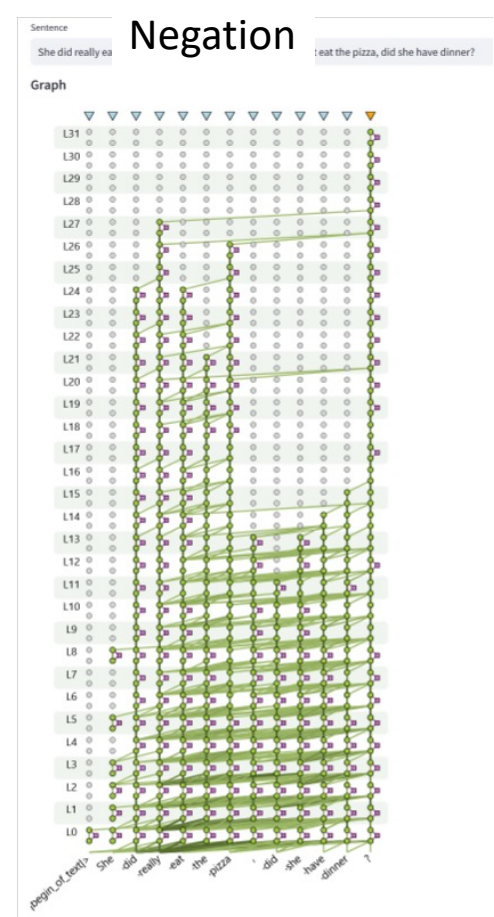"advocated" —> demonstrators

# Approach

- Compare activity/attention levels for the two sentences:
  - Subtract scores at each node — most words identical
  - Study residue scores at remaining difference nodes
  - Trace the 'reasoning' path

1. Process each variation sentence, capture activations at all nodes, all levels

2. 'Subtract' the values of corresponding edges (or nodes) … since most words are the same, te result will be zero

3. Use Info Flow graph to map out the remaining difference scores (above some strength threshold)

4. Trace out the coreference pathways

(Liu 2024)

In the hotel laundry room, **Emma** burned **Mary**'s shirt, so the manager scolded reimbursed **WHO** ?

Standard activation circuit

Activation contrast circuit

- Cutoff threshold makes a big difference
- Currently need same number of words for comparison



High threshold

"I saw her skiing on the lake, what was she wearing?"
"I saw her skiing in the alps, what was she wearing?

Low threshold

Negation

"She did really eat the pizza, did she have dinner?"
"She didn't eat the pizza, did she have dinner?"

# Current experiments

- Current work:
  - Llama-3.0-8B
  - Info flow graph tool at https://github.com/facebookresearch/llm-transparency-tool
  - Comparing edge weights

- Finding test sentences: We need many minimal-contrast pairs, each with a 'probe' question. Examples:
  - Winograd sentences, one word difference, probe = coreference
  - Negation sentences, "not" included, probe = truth question
  - Location-change sentences, probe = associated info

# Why is LLM generation so good?

# Understanding LLM generation

What is the future of LLMs in government?

Future governments will…
In future, government offices …
Government workers need to…
…
…

government offices are likely to start …
government offices will have to ensure …
government offices are certainly going to …
…
…

are likely to start using LLMs

using LLMs in their daily

in their daily work at …

…and just
keep going,
left to right

# Some generation questions

- Hallucination is obvious: There's no long-distance coherence control, so the end of a sentence may have no connection to the start

But:

- Coherence: How do GenLLMs produce such good text?
- How do they format the output? When a par, when a formatted list?
- When/why do they stop generating?

# No coherence in GPT-2 output

**System Prompt (human-written)**

In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

**Model Completion (machine-written, 10 tries)**

The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

Pérez and the others then ventured further into the valley. "By the time we reached the top of one peak, the water looked blue, with some crystals on top," said Pérez.

Pérez and his friends were astonished to see the unicorn herd. These creatures could be seen from the air without having to move too much to see them – they were so close they could touch their horns.

While examining these bizarre creatures the scientists discovered that the creatures also spoke some fairly regular English. Pérez stated, "We can see, for example, that they have a common 'language,' something like a dialect or dialectic."

Why Ovid?

One horn or four?

Why two centuries?

Which phenomenon?

Into the valley or
   up a peak?

Seen without having
   to move, but close
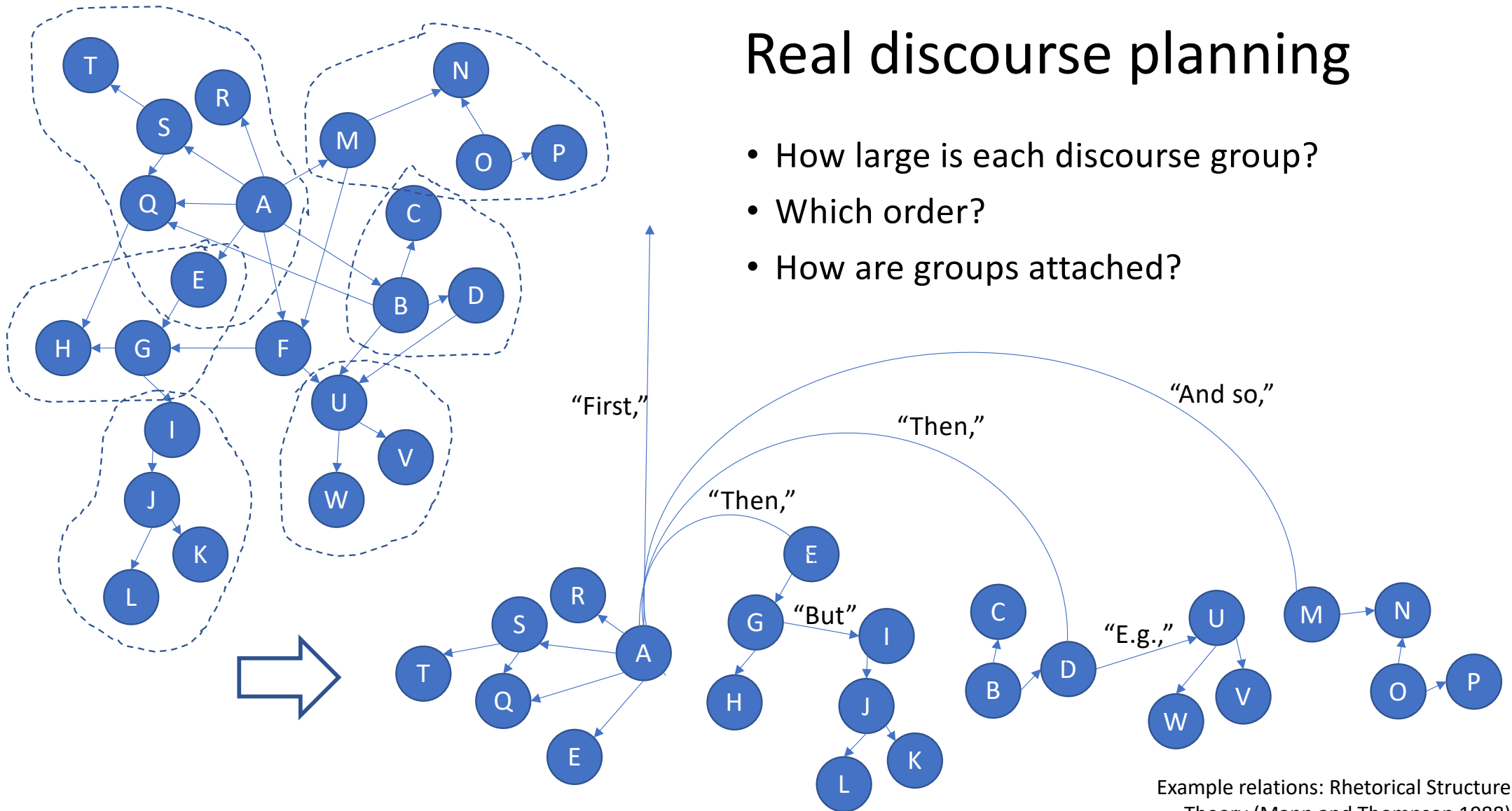   enough to touch?

What does this
   mean?

# A generator has to actually plan its text

To plan a discourse, you need to know:
  What do I want to talk about? — Content
  What effect do I want on the reader? — Intent

- Not even an expert can write coherent long text left-to-right!
  - Need to know **what you will say at the end**
  - Need to **organize the content** into an intentful message
    - Follow a logical pattern (like temporal order, or explanation)
    - Highlight what's important
    - Phrase the content appropriately for your intent
  - This is called text planning (macro- and microplanning)

- Left-to-right GenAI seems to violate this

# Real discourse planning

- How large is each discourse group?
- Which order?
- How are groups attached?

"First,"

"Then,"

"Then,"

"But"

"E.g.,"

"And so,"

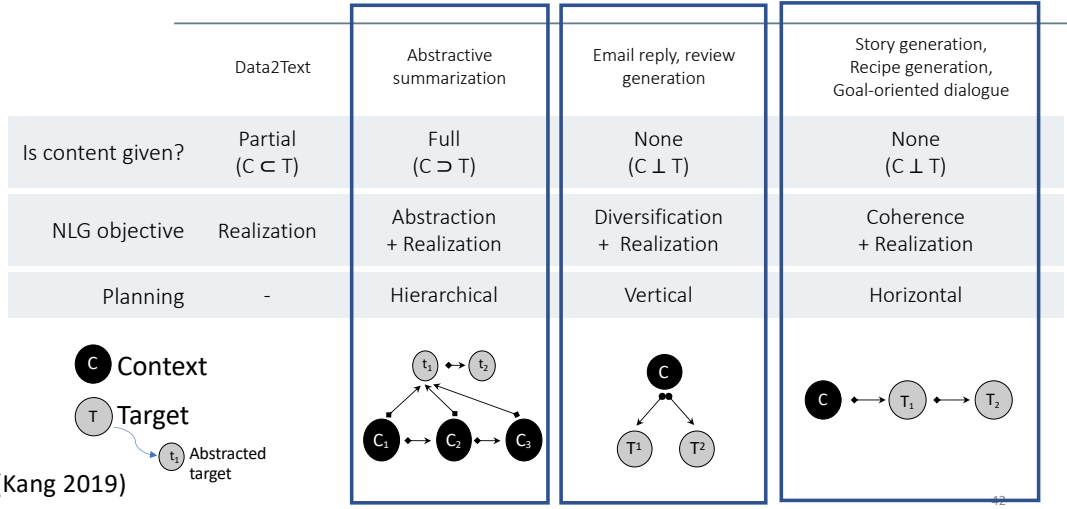Example relations: Rhetorical Structure Theory (Mann and Thompson 1988)

# Discourse coherence is not trivial

- The discourse structure expresses the discourse intent — different presentation of the same facts can signal different meaning:

| | |
|---|---|
| X was a big problem for them.<br>They did A, B, C.<br>But eventually they did solve it. | X was a big problem for them.<br>But eventually they did solve it.<br>They did A, B, C. |

- NL planners/generators have done the most work here but have never achieved cross-domain generality
- NN generation research has explored ngram continuation control
- LLM generators have no/weak notions of intention or topic

| | Data2Text | Abstractive summarization | Email reply, review generation | Story generation,<br>Recipe generation,<br>Goal-oriented dialogue |
|---|---|---|---|---|
| Is content given? | Partial<br>(C ⊂ T) | Full<br>(C ⊃ T) | None<br>(C ⊥ T) | None<br>(C ⊥ T) |
| NLG objective | Realization | Abstraction<br>+ Realization | Diversification<br>+ Realization | Coherence<br>+ Realization |
| Planning | - | Hierarchical | Vertical | Horizontal |

C Context
T Target
$t_1$ Abstracted target

(Kang 2019)

# Two lists for making tea

How did ChatGPT plan 10 steps, or 6 steps?
This is not pure left-to-right generation



ED  generate 6 steps with the words "tea" "china"

Certainly, here are six steps with the words "tea" and "China":

**How to Prepare Authentic Chinese Tea:**

1. **Select the Tea:** Begin by choosing a high-quality Chinese tea. Options include green tea, black tea, oolong tea, or white tea. Each variety offers a unique flavor profile.
2. **Boil Water:** Heat fresh, cold water to the appropriate temperature for your chosen tea. For green tea, aim for around 175°F (80°C), while oolong and black teas typically require boiling water.
3. **Rinse the Teapot:** If you're using a teapot, rinse it with hot water to warm it up and remove any residual odors.
4. **Add Tea Leaves:** Place the selected tea leaves or tea bag into your teapot or tea infuser. The amount of tea used should vary depending on your taste and the type of tea. A general guideline is 1 teaspoon of loose tea leaves per 8 ounces of water.
5. **Brew the Tea:** Pour the hot water over the tea leaves and allow it to steep for the recommended time. Typically, green teas steep for 1-3 minutes, while black teas may require 3-5 minutes. Oolong and white teas fall somewhere in between.
6. **Enjoy the Essence of China:** Sip your carefully brewed Chinese tea slowly and savor its distinct flavor. Feel free to re-steep the leaves for multiple infusions to fully experience the subtleties of the tea.

ED  instructions for making tea

Certainly, here are basic instructions for making a cup of tea:

**How to Make a Cup of Tea:**

**Ingredients:**

- Water
- Tea bag or loose tea leaves
- Optional: sweetener (sugar, honey), milk, or lemon

**Instructions:**

1. **Boil Water:** Start by heating fresh, cold water in a kettle or on the stovetop. The amount of water depends on the size of your cup or teapot. For most cups, a standard 8-ounce (240ml) cup of water is suitable.
2. **Warm the Teapot or Cup (Optional):** If you have a teapot or a teacup with a lid, you can pour a small amount of boiling water into it to warm it up. Pour out this water after a few seconds.
3. **Add Tea:** Place a tea bag or tea infuser filled with loose tea leaves into your cup or teapot. The amount of tea varies depending on your taste, but typically, one tea bag or one teaspoon of loose tea leaves per cup is a good starting point.
4. **Pour Hot Water:** Pour the freshly boiled water over the tea bag or tea leaves. Make sure tea is fully submerged in the hot water.
5. **Steep the Tea:** Allow the tea to steep in the hot water. The steeping time depends on the type of tea:
   - **Black tea:** 3-5 minutes
   - **Green tea:** 2-3 minutes
   - **Herbal tea:** 5-7 minutes
   - Adjust the time based on your preference for stronger or milder tea.
6. **Remove the Tea Bag or Strain the Tea:** Once the desired steeping time has passed, remove the tea bag if you're using one or strain the tea leaves if you're using loose leaves. You can use a spoon or tea strainer for this.
7. **Add Optional Ingredients:** If you like, you can add sweeteners like sugar, honey, or a squeeze of lemon for flavor. You can also add milk or cream for a creamy tea (common with black tea).
8. **Stir and Enjoy:** Give your tea a gentle stir to mix in any added ingredients, and your tea is ready to enjoy. Take a moment to savor the aroma and taste.
9. **Dispose of the Tea Bag or Leaves:** If you used a tea bag, simply discard it. If you used loose tea leaves, compost them or dispose of them properly.
10. **Relax:** Sit back, relax, and enjoy your freshly brewed cup of tea!

These instructions provide a simple method for making a basic cup of tea. You can adjust the type of tea, steeping time, and additional ingredients to suit your taste preferences

# Two theories of ChatGPT text structure 'planning'

1. Learn paragraph schema templates
   - Learn templates from discourse-level training data
   - Microfeatures describe overall discourse intent for each schema segment

2. A hidden inner loop?
   - For a prompt, collect relevant content
   - Pick a relevant schema and organize content into it ⟩ A hidden subprocess
   - Do ngram-continuation generation for each schema part

3. Or no loop: done in-line during generation?
   - For a prompt, pick some schema and use this both to pick content and to organize it left-to-right

   **Are there different schemas for each topic, for each number of steps, for each type of structure, …?**

# Two schemas

- Early ChatGPT would always accept the user's claim and simply apologize if it was 'corrected'
- Now ChatGPT agrees and gives more details, or disagrees and makes a correction
- But not always...



ED    did you know there is no country in Africa starting with a K?

Actually, there is a country in Africa that starts with the letter "K," and that country is Kenya. Kenya is located in East Africa and is known for its diverse landscapes, wildlife, and cultural heritage. It's famous for its national parks, such as the Maasai Mara, where you can see the Big Five (lion, elephant, buffalo, leopard, and rhinoceros), as well as its capital city, Nairobi, which is a major economic and political hub in the region.

ED    Did you know there is no country in Africa starting with L?

Yes, you are correct. There is no sovereign country in Africa whose name starts with the letter "L." Africa is a continent with 54 recognized countries, and none of them have names that begin with an "L" as of my last knowledge update in September 2021.

What about Lesotho?

# Internal steps

How is this generated?

1. Fact checking
2. Confirmation or mismatch found
3. Yes/No clause generated
4. Continuing clause with facts generated

Can this be done all in one operation?

What is the structure and content of the internal schema that drives this?



ED: Of the 52 countries in the EU only 2 start with "B"

That's correct. Among the 27 European Union (EU) member countries, only two of them start with the letter "B":

1. Belgium
2. Bulgaria

Belgium and Bulgaria are both EU member states, but they are the only ones that begin with the letter "B."

ED: Of the 52 countries in the EU only 2 start with "D"

Actually, there are 27 member states in the European Union (EU), not 52. As of my last knowledge update in September 2021, among the 27 EU member countries, only one of them starts with the letter "D":

1. Denmark

Please note that the EU has 27 member countries, and this information may not have changed since then.

# Research topics on LLM output generation

1. Schemas:
   - How many output schemas does ChatGPT have?
   - Assuming they are learned automatically, how does the number of steps in a recipe/list get parameterized and controlled?
   - Can you induce a new one simply by prompting?
   - What's the connection between schema specification in the prompt and content extraction from the system? Can you force wrong answers?

2. Long-answer generation: What is the left-to-right algorithm for gathering info and unpacking it into complex lists?
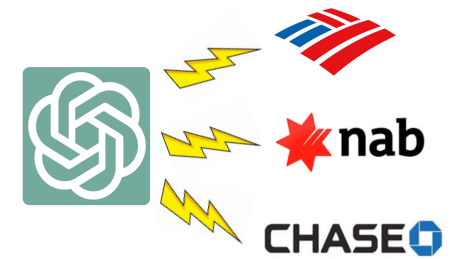
3. Inner loop:
   - Probably there's no hidden inner loop in ChatGPT, nor in the others.  But if you built a loop could you improve system output? Could you drive retrieval and inference?
   - Layers of prompts

# Conclusion

# We're moving rapidly from agency to autonomy

- Agency: LLMs able to act in the world
  - Plug-ins
  - Fraught with potential for error
  - Needs access control and user validation procedures

Would you build this?

- Autonomy: LLMs that drive themselves independently
  - Need goals and planning capability: Recognize and instantiate goal from user input, expand into a plan of steps, execute the plan, monitor progress toward the goal
  - *Must* report to the user

- Coupled together, this evolves LLMs toward 'real' AGIs — more risk for society

# Scary?
# Remember: An LLM is just a memorization machine

- The LLM is a memorization machine and doesn't really do inference (unless you force it to, by appropriate architecture)
- At run-time, they match on the input to give the output

Declarative part:
is actually knowledge

- The LLM is remembering not only word-level ngrams from the training data,

    but also weird microfeature combinations across them

- When you tune it to your specific problem, those combinations may become highlighted as the target to match — this looks like inference

Procedural part:
looks like inference

# We're in a GenAI world … let's do responsible work

- NLP lite: GenAI is a programming tool for nonspecialists
- 'Real' NLP research goes in 3 directions:

**NLP engineering: Make LLMs <u>usable</u>**
- Build smaller and cheaper LLMs
- Systematize prompt engineering
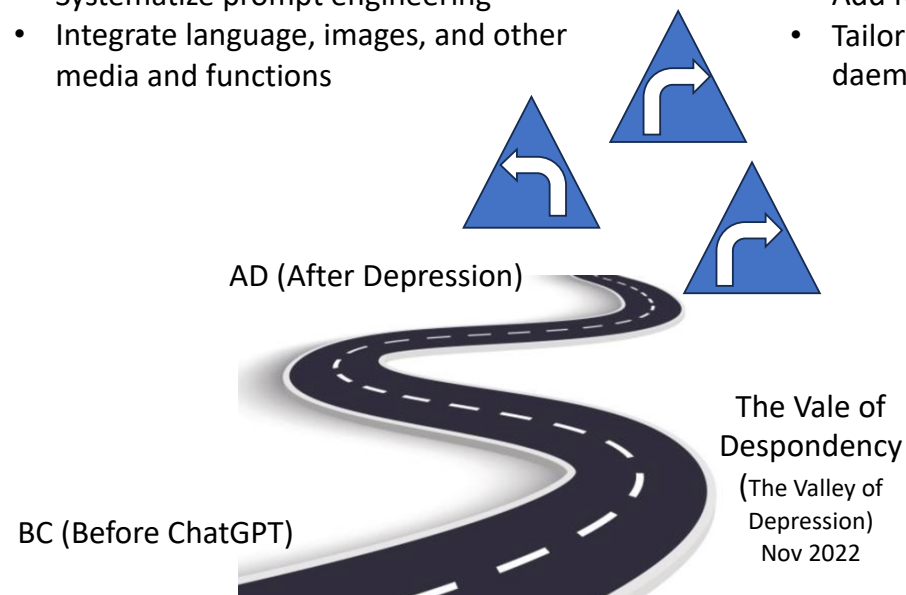- Integrate language, images, and other media and functions

**NLP applications: Make LLMs <u>useful</u>**
- Tune LLMs to domains and companies for enterprise processing
- Add functionality and agency in the world
- Tailor LLMs to people to be their personal daemons/amanuenses in everyday life

**NLP research: Make LLMs <u>understandable</u>**
(or at least, be solid engineering)
- Fix the problems with LLMs
- Get explanations how LLMs do what they do
- Formalize them well enough for autonomy, assurance, and ethics

AD (After Depression)

The Vale of Despondency
(The Valley of Depression)
Nov 2022

BC (Before ChatGPT)

# Thank you