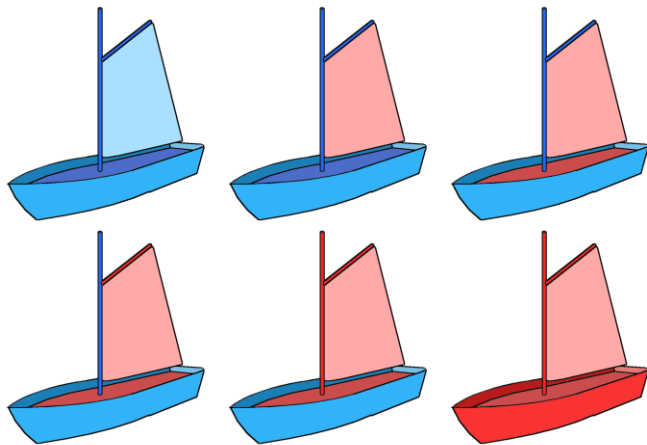# Discourse analysis as a diagnostic lens:
*Untangling some of the riddles complicating LLM evaluations*

**SEPTEMBER, 2025**

Samantha Finkelstein, PhD
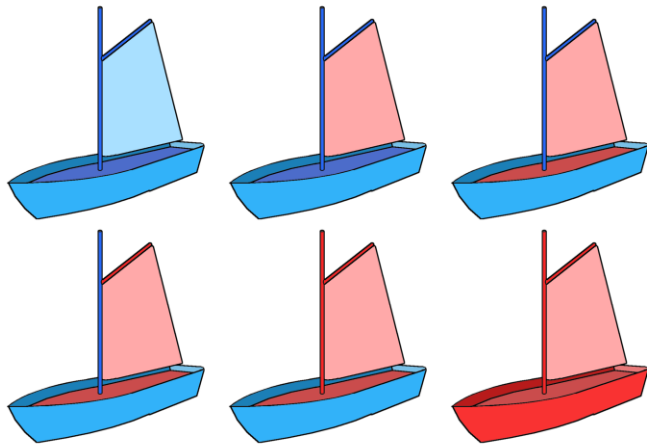Sr. Human-Centered AI Research Scientist,

Advancing Software for National Security

# What are we trying to evaluate?

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

2

# What are we trying to evaluate?

*The Ship of Theseus is a philosophy paradox:*
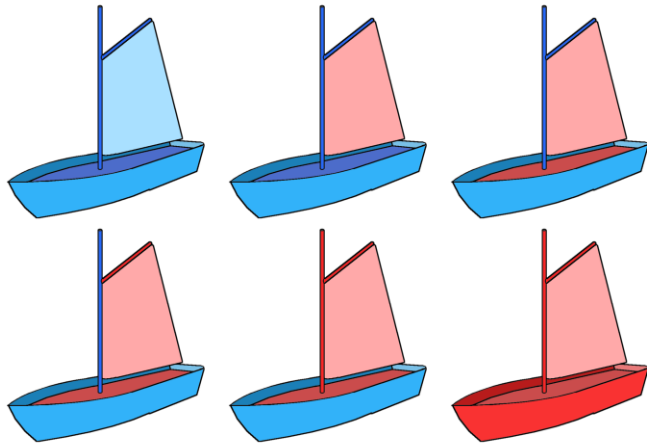


*If we replace one plank in the ship of Theseus, is it the same ship?*

*What if we've replaced every plank?*

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

Advancing Software for National Security

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

3

# What are we trying to evaluate?

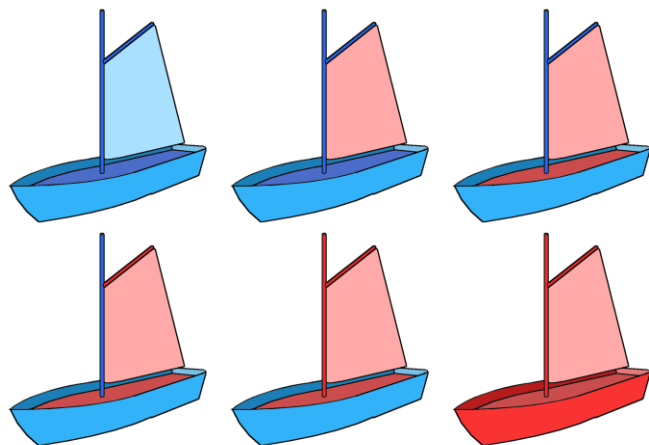*The Ship of Theseus is a philosophy paradox:*



*If we replace one plank in the ship of Theseus, is it the same ship?*

*What if we've replaced every plank?*

- If "**same**" means *the atoms that make up its structure,* then *no*
- If "**same**" means *the functionality that it provides,* then *yes*

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

4

# What are we trying to evaluate?

*The Ship of Theseus is a philosophy paradox:*

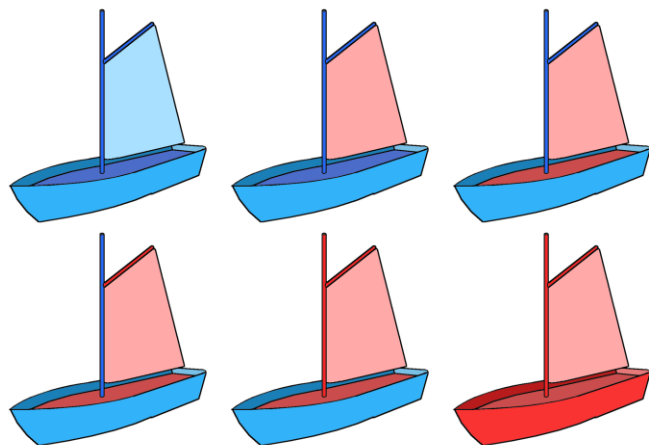*If we replace one plank in the ship of Theseus, is it the same ship?*

*What if we've replaced every plank?*

- If "**same**" means *the atoms that make up its structure,* then *no*
- If "**same**" means *the functionality that it provides,* then *yes*

**The wording of the question obscures the very distinctions you most need to center!**

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

5

# What are we trying to evaluate?

*The Ship of Theseus is a* ~~philosophy paradox~~ *definition trick:*



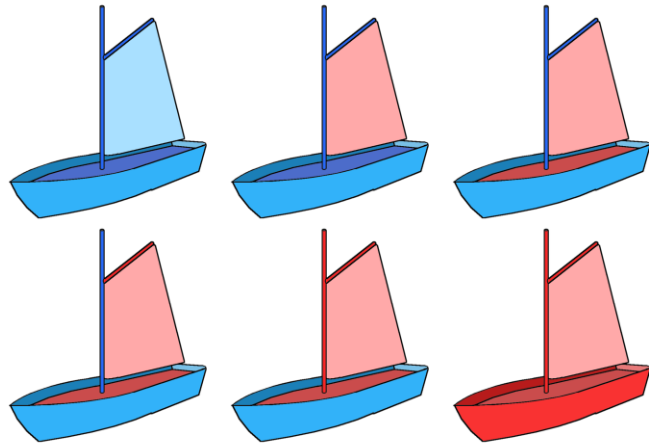*If we replace one plank in the ship of Theseus, is it the same ship?*

*What if we've replaced every plank?*

- If "**same**" means *the atoms that make up its structure,* then *no*
- If "**same**" means *the functionality that it provides,* then *yes*

**The wording of the question obscures the very distinctions you most need to center!**

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

6

# What are we trying to evaluate?

*Many LLM Evaluation "riddles" are similarly a definition trick:*



*If we replace one plank in the ship of Theseus, is it the same ship?*
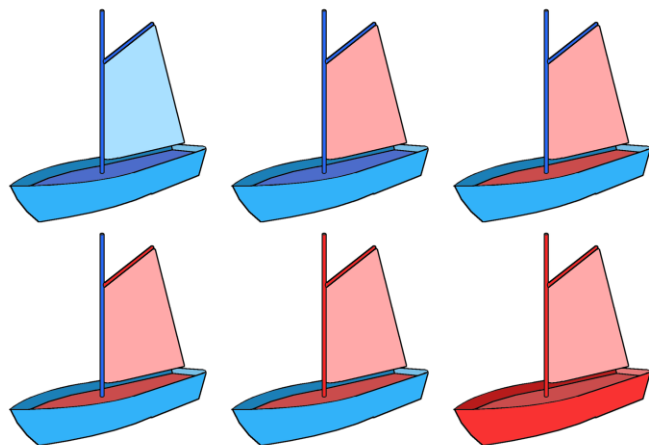
*What if we've replaced every plank?*

- If "**same**" means *the atoms that make up its structure,* then *no*
- If "**same**" means *the functionality that it provides,* then *yes*

**The wording of the question obscures the very distinctions you most need to center!**

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

7

# What are we trying to evaluate?

*Many LLM Evaluation "riddles" are similarly a definition trick:*



*Are LLMs trustworthy?*

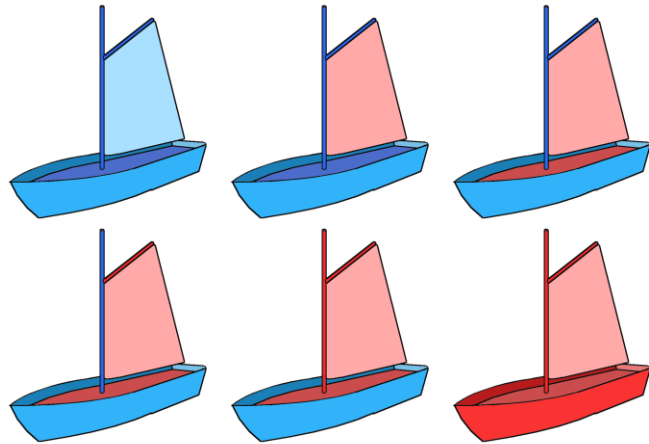*Are they reliable? Are they safe?*

- If "**same**" means *the atoms that make up its structure,* then *no*
- If "**same**" means *the functionality that it provides,* then *yes*

> **The wording of the question obscures the very distinctions you most need to center!**

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

8

# What are we trying to evaluate?

*Many LLM Evaluation "riddles" are similarly a definition trick:*



*Are LLMs trustworthy?*

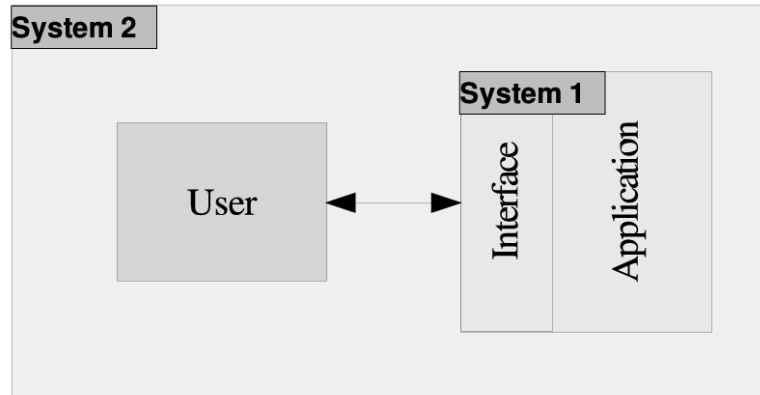**What do we mean by trustworthy?**

*Are they reliable? Are they safe?*

**What do we mean by reliable? By safe?**

- If "**same**" means *the atoms that make up its structure,* then *no*
- If "**same**" means *the functionality that it provides,* then *yes*

**The wording of the question obscures the very distinctions you most need to center!**

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

9

# What are we trying to evaluate?

*Many LLM Evaluation "riddles" are similarly a definition trick:*



*Are LLMs trustworthy?*

**What do we mean by trustworthy?**
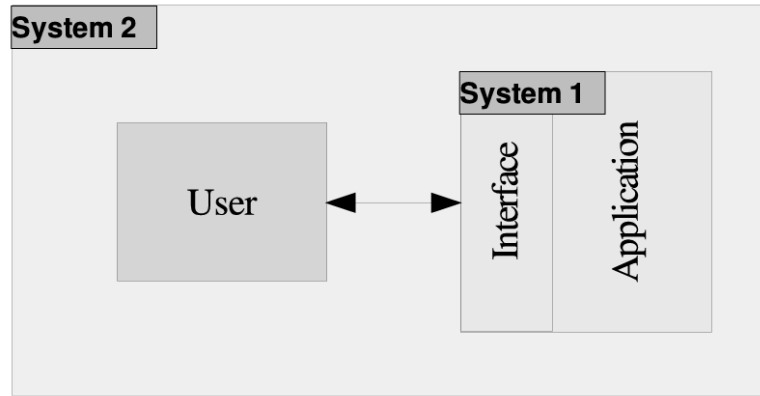
*Are they reliable? Are they safe?*

**What do we mean by reliable? By safe?**

- If "**same**" means *the atoms that make up its structure,* then *no*
- If "**same**" means *the functionality that it provides,* then *yes*

**The wording of the question obscures the very distinctions you most need to center!**

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

10

# What are we trying to evaluate?

*Many LLM Evaluation "riddles" are similarly a definition trick:*



### Are LLMs trustworthy?
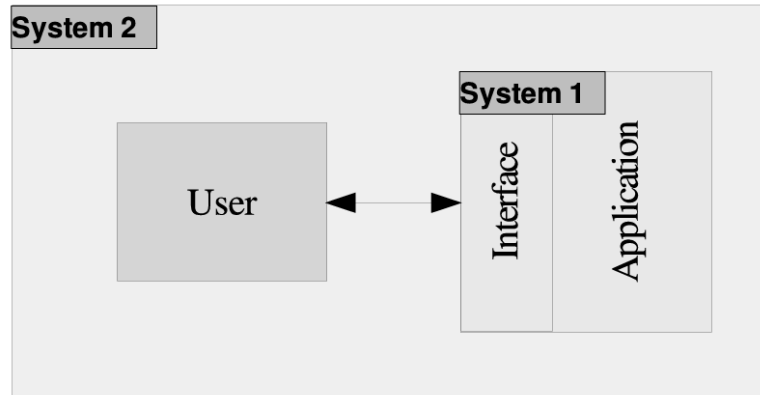
**What do we mean by trustworthy?**

### Are they reliable? Are they safe?

**What do we mean by reliable? By safe?**

- Are you evaluating **"system performance"** at the level of its components (decontextualized)?
- Or **"system performance"** at the level of its *impacts* (contextualized)?

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

11

# What are we trying to evaluate?

*Many LLM Evaluation "riddles" are similarly a definition trick:*



## System 1: LLM capabilities
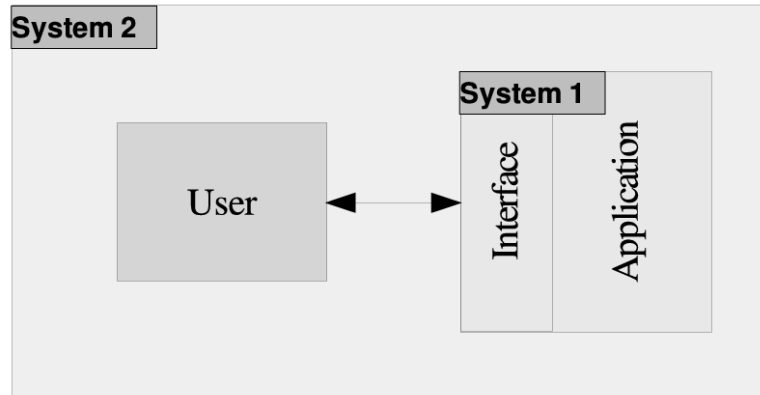**What are the quality attributes of this _language model_?**

## System 2: Application functionalities
**What are the quality attributes of this _application_** (which is powered, in part, by a language model)?

- Are you evaluating **"system performance"** at the level of its components (decontextualized)?
- Or **"system performance"** at the level of its *impacts* (contextualized)?

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

12

# What are we trying to evaluate?

*Many LLM Evaluation "riddles" are similarly a definition trick:*



## System 1: LLM capabilities
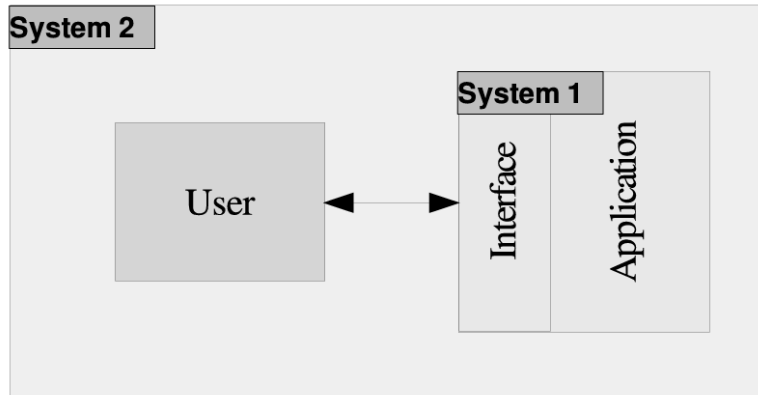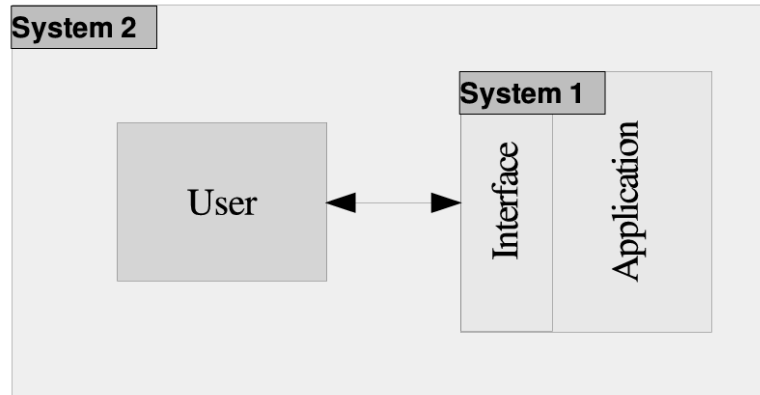**Processing speed! Rouge score! Token weights!**

## System 2: Application functionalities
**Benefits! Risks! Effectiveness! Usefulness! Safety!**

- Are you evaluating **"system performance"** at the level of its components (decontextualized)?
- Or **"system performance"** at the level of its *impacts* (contextualized)?

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

13

# What are we trying to evaluate?

*System 2 evaluation is necessary for operational deployment decisions*



System 1: LLM capabilities
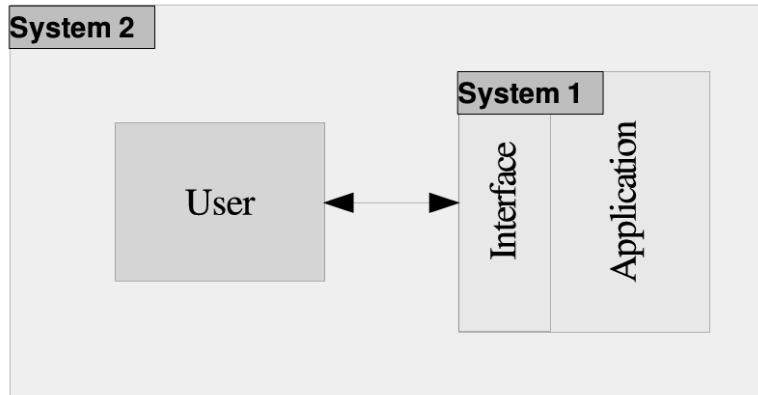
*Processing speed! Rouge score! Token weights!*

*System 2: Application functionalities*

*Benefits! Risks! Effectiveness! Usefulness! Safety!*

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

14

# What are we trying to evaluate?

## System 2 evaluation is necessary for operational deployment decisions



*System 1: LLM capabilities*

*Processing speed! Rouge score! Token weights!*

*System 2: Application functionalities*
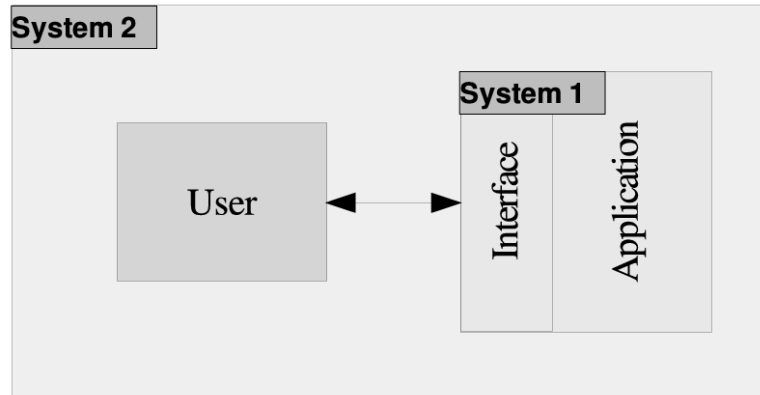
*Benefits! Risks! Effectiveness! Usefulness! Safety!*

*The rest of this talk is designed to make it easier for you to "solve the riddle" of what this means for your own deployments*

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

15

# What are we trying to evaluate?

## System 2 evaluation is necessary for operational deployment decisions



*System 1: LLM capabilities*
*Processing speed! Rouge score! Token weights!*

*System 2: Application functionalities*
**Benefits! Risks! Effectiveness! Usefulness! Safety!**

*The rest of this talk is designed to make it easier for you to "solve the riddle" of what this means for your own deployments*

1. Define the menu of these functionalities
2. Share evaluation considerations for each

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

16

# What are we trying to evaluate?

## System 2 evaluation is necessary for operational deployment decisions



System 1: LLM capabilities
Processing speed! Rouge score! Token weights!

## System 2: Application functionalities
Benefits! Risks! Effectiveness! Usefulness! Safety!

*The rest of this talk is designed to make it easier for you to "solve the riddle" of what this means for your own deployments*

1. Define the menu of these functionalities
2. Share evaluation considerations for each
3. Convince you about discourse analysis ☺

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

17

# LLMs provide 3 categories of application functionality

## Conversation

System enables a dialogic interaction where users construct input - and interpret output - through the lens of discourse.
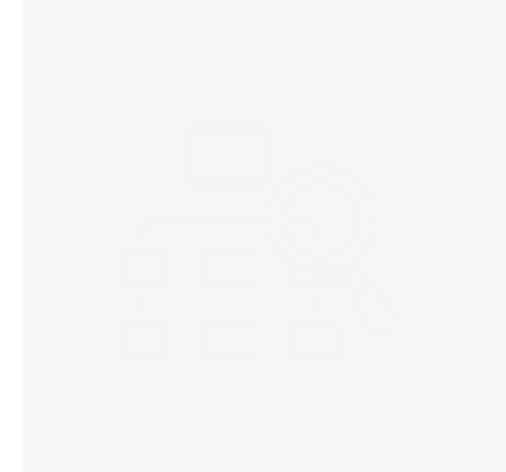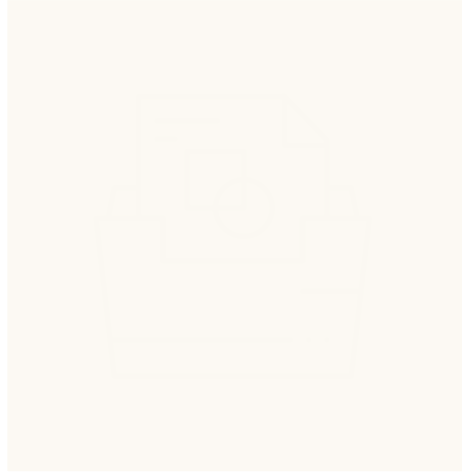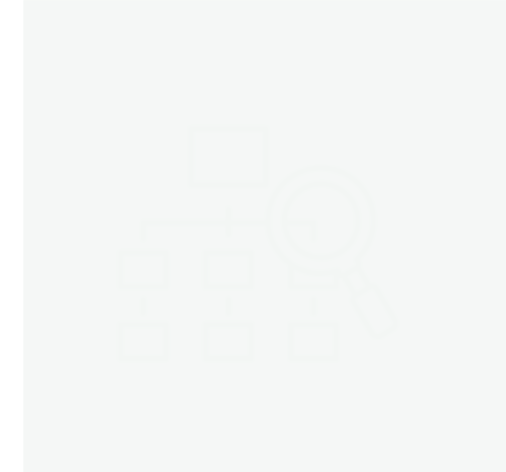
## Generation

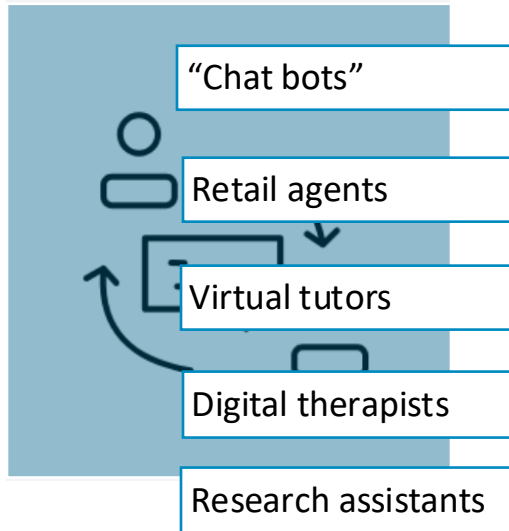System enables user specification of criteria on which to deliver a stand-alone artifact.

## Analysis

System enables transformation of language signals into different signals, in accordance with identified specifications.

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

18

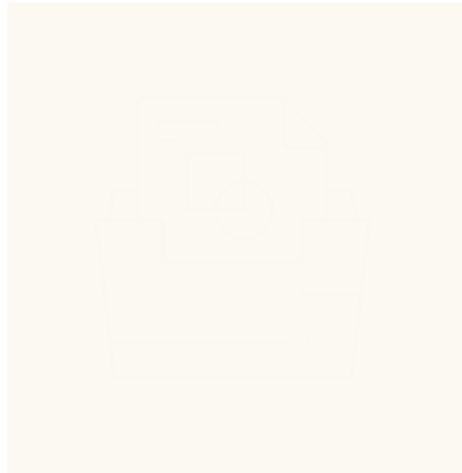# LLMs provide 3 categories of application functionality

## Conversation

*System enables a dialogic interaction where users construct input - and interpret output - through the lens of discourse.*
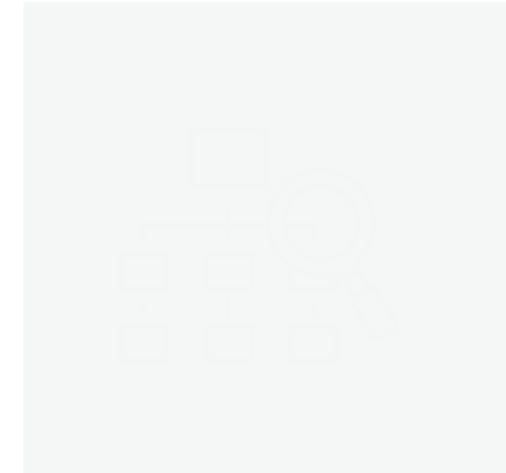
## Generation

*System enables user specification of criteria on which to deliver a stand-alone artifact.*
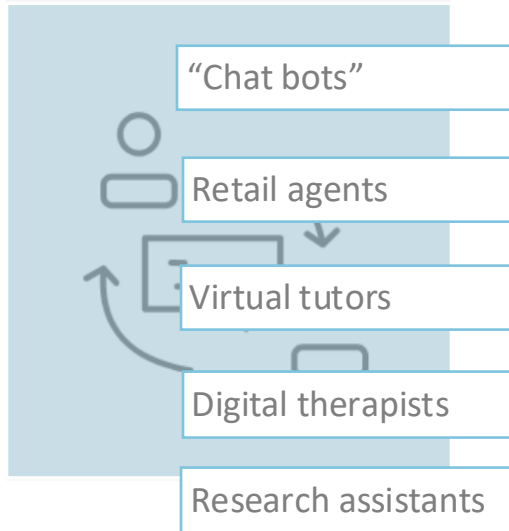
## Analysis

*System enables transformation of language signals into different signals, in accordance with identified specifications.*

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

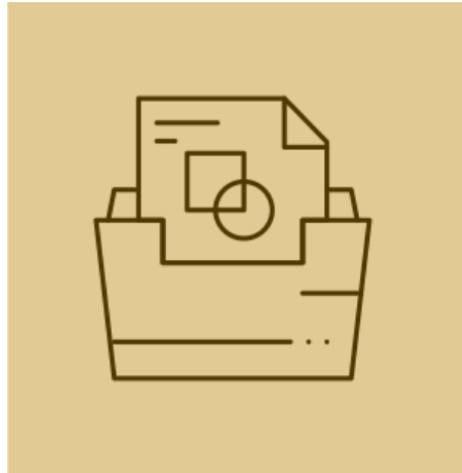19

# LLMs provide 3 categories of application functionality

## Conversation

*System enables a dialogic interaction where users construct input - and interpret output - through the lens of discourse.*

"Chat bots"

Retail agents

Virtual tutors

Digital therapists

Research assistants

## Generation

*System enables user specification of criteria on which to deliver a stand-alone artifact.*

## Analysis

*System enables transformation of language signals into different signals, in accordance with identified specifications.*

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

20

# LLMs provide 3 categories of application functionality



## Conversation

*System enables a dialogic interaction where users construct input - and interpret output - through the lens of discourse.*

- "Chat bots"
- Retail agents
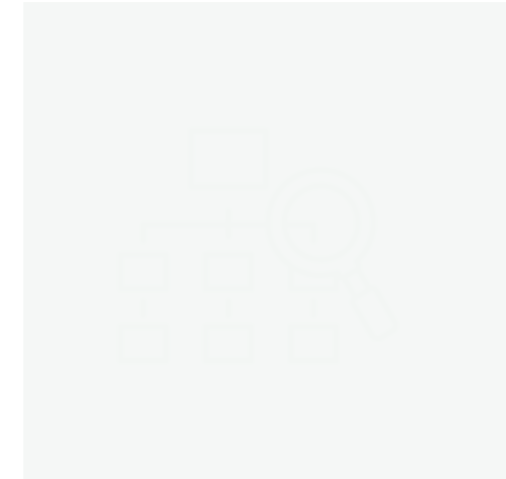- Virtual tutors
- Digital therapists
- Research assistants

## Generation

*System enables user specification of criteria on which to deliver a stand-alone artifact.*

## Analysis
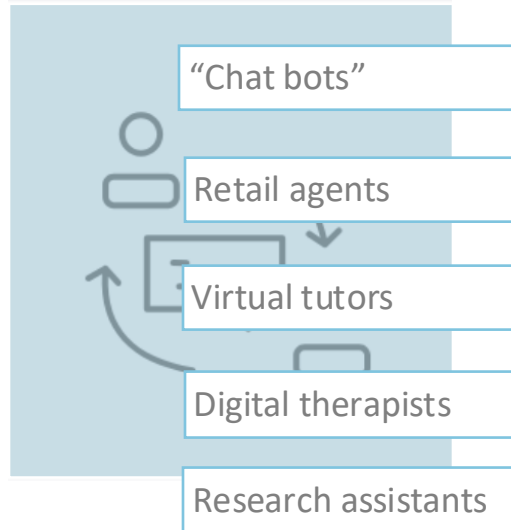
*System enables transformation of language signals into different signals, in accordance with identified specifications.*

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

21
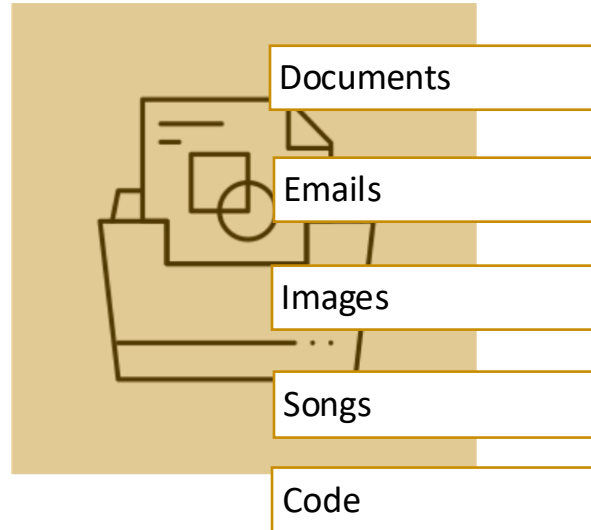
# LLMs provide 3 categories of application functionality

## Conversation

*System enables a dialogic interaction where users construct input - and interpret output - through the lens of discourse.*
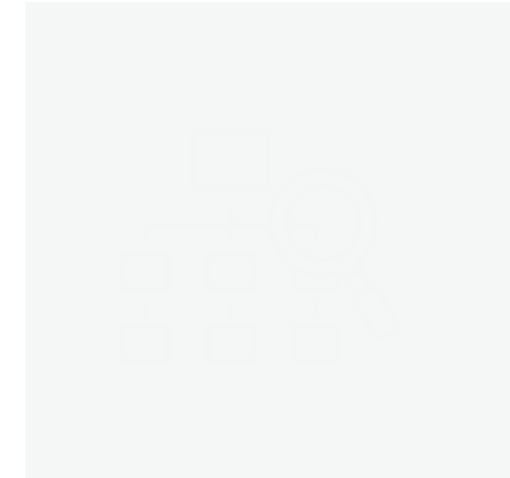
- "Chat bots"
- Retail agents
- Virtual tutors
- Digital therapists
- Research assistants

## Generation

*System enables user specification of criteria on which to deliver a stand-alone artifact.*
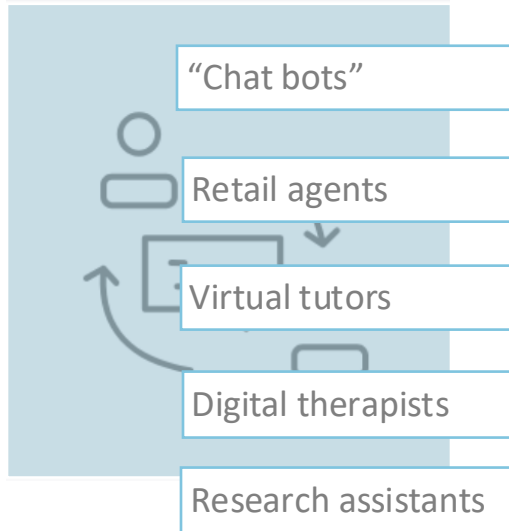
- Documents
- Emails
- Images
- Songs
- Code

## Analysis

*System enables transformation of language signals into different signals, in accordance with identified specifications.*

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

22
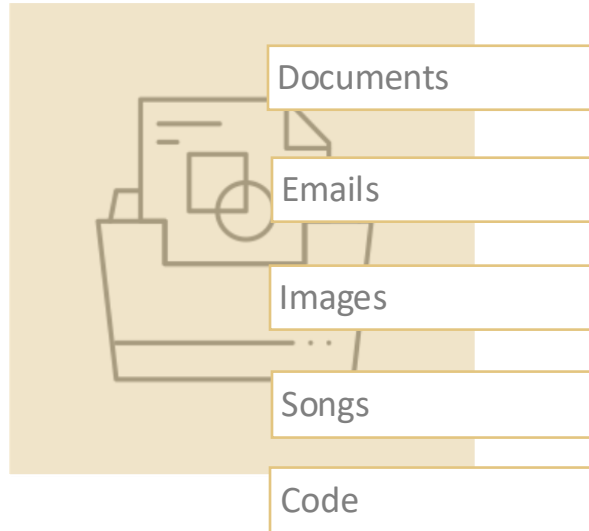
# LLMs provide 3 categories of application functionality

## Conversation

*System enables a dialogic interaction where users construct input - and interpret output - through the lens of discourse.*
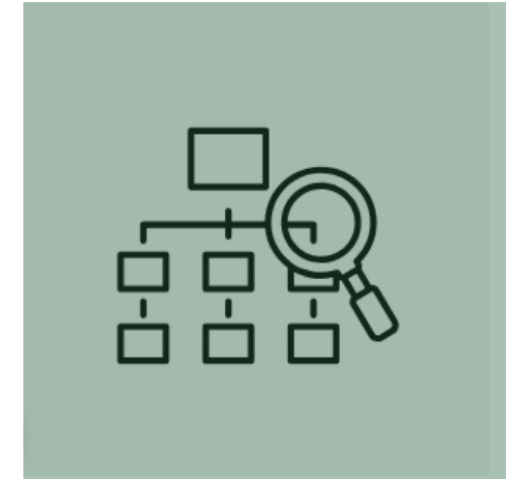
- "Chat bots"
- Retail agents
- Virtual tutors
- Digital therapists
- Research assistants

## Generation

*System enables user specification of criteria on which to deliver a stand-alone artifact.*
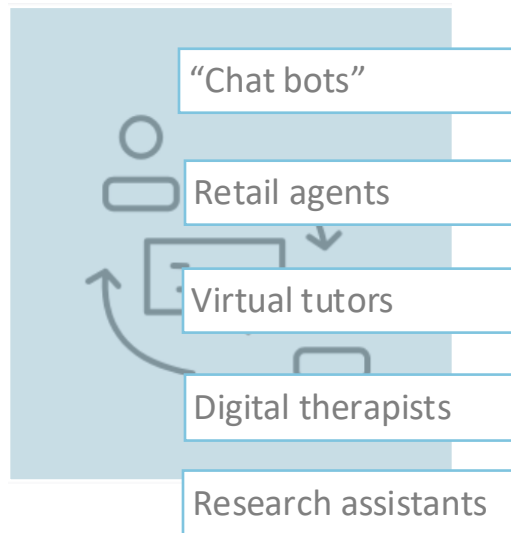
- Documents
- Emails
- Images
- Songs
- Code

## Analysis

*System enables transformation of language signals into different signals, in accordance with identified specifications.*

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

23
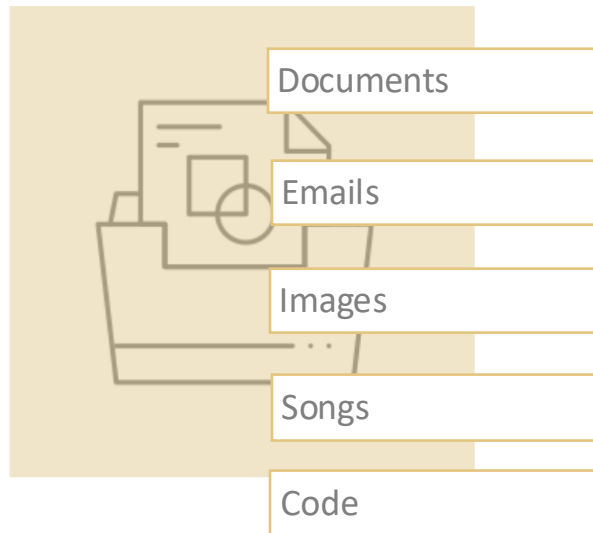
# LLMs provide 3 categories of application functionality

## Conversation

*System enables a dialogic interaction where users construct input - and interpret output - through the lens of discourse.*
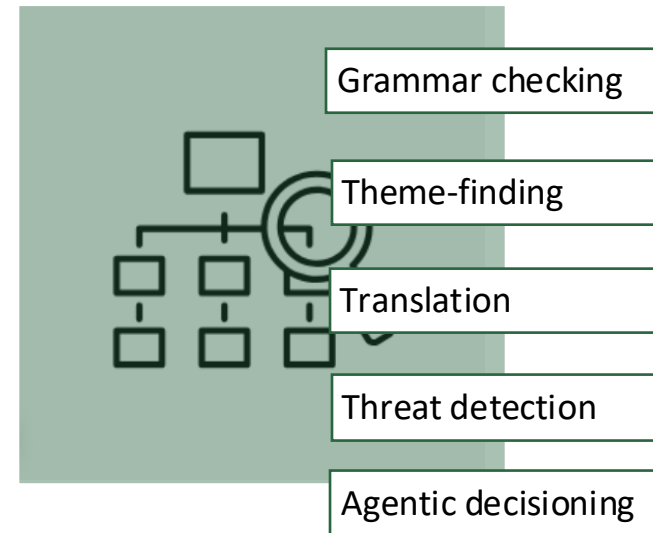
- "Chat bots"
- Retail agents
- Virtual tutors
- Digital therapists
- Research assistants

## Generation

*System enables user specification of criteria on which to deliver a stand-alone artifact.*

- Documents
- Emails
- Images
- Songs
- Code

## Analysis

*System enables transformation of language signals into different signals, in accordance with identified specifications.*

- Grammar checking
- Theme-finding
- Translation
- Threat detection
- Agentic decisioning

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

24

# Let's apply this to an example: intelligence analysis

## Conversation

**Intelligence ideation**

The analyst and LLM discuss intel across sources to strengthen the interpretive scope and rigor.
-------

## Generation

**Document summarization**

The analyst feeds in a long intelligence report (or set of reports), and the LLM generates a summary that retains the "most important" information, in accordance with specifications.
-------

## Analysis

**Propaganda detection**

The system ingests streams of data (e.g., sourced reports, messages, news articles) and uses an LLM to examine those documents for signals of potential adversarial propaganda or influence.
-------

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

25

# Let's apply this to an example: intelligence analysis

## Conversation

### Intelligence ideation

The analyst and LLM discuss intel across sources to strengthen the interpretive scope and rigor.

-------

The LLM can support structure and ideation across considerations like:
- **Brainstorming** (identify and interrogate alternative explanations)
- **Surfacing ambiguities** (identifying blind spots, open questions, testing edge-cases, poking at assumptions)
- **Contextualizing** (applying intel across situations to reveal subtle patterns or applications)

## Generation

### Document summarization

The analyst feeds in a long intelligence report (or set of reports), and the LLM generates a summary that retains the "most important" information, in accordance with specifications.

-------

## Analysis

### Propaganda detection

The system ingests streams of data (e.g., sourced reports, messages, news articles) and uses an LLM to examine those documents for signals of potential adversarial propaganda or influence.

-------

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

26

# Let's apply this to an example: intelligence analysis

## Conversation

**Intelligence ideation**

The analyst and LLM discuss intel across sources to strengthen the interpretive scope and rigor.

------

The LLM can support structure and ideation across considerations like:

- **Brainstorming** *(identify and interrogate alternative explanations)*
- **Surfacing ambiguities** *(identifying blind spots, open questions, testing edge-cases, poking at assumptions)*
- **Contextualizing** *(applying intel across situations to reveal subtle patterns or applications)*

## Generation

**Document summarization**

The analyst feeds in a long intelligence report (or set of reports), and the LLM generates a summary that retains the "most important" information, in accordance with specifications.

------

Those specifications can span:

- **importance** *(scrutiny / prioritization appropriate across detail types)*
- **phrasing** *(the level of paraphrasing allowable / desirable)*
- **formatting** *(output requirements)*

## Analysis

**Propaganda detection**

The system ingests streams of data (e.g., sourced reports, messages, news articles) and uses an LLM to examine those documents for signals of potential adversarial propaganda or influence.

------

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

27

# Let's apply this to an example: intelligence analysis

## Conaversation

**Intelligence ideation**

The analyst and LLM discuss intel across sources to strengthen the interpretive scope and rigor.
-------

The LLM can support structure and ideation across considerations like:
- **Brainstorming** *(identify and interrogate alternative explanations)*
- **Surfacing ambiguities** *(identifying blind spots, open questions, testing edge-cases, poking at assumptions)*
- **Contextualizing** *(applying intel across situations to reveal subtle patterns or applications)*

## Generation

**Document summarization**

The analyst feeds in a long intelligence report (or set of reports), and the LLM generates a summary that retains the "most important" information, in accordance with specifications.
------

Those specifications can span:
- **importance** *(scrutiny / prioritization appropriate across detail types)*
- **phrasing** *(the level of paraphrasing allowable / desirable)*
- **formatting** *(output requirements)*

## Analysis

**Propaganda detection**

The system ingests streams of data (e.g., sourced reports, messages, news articles) and uses an LLM to examine those documents for signals of potential adversarial propaganda or influence.
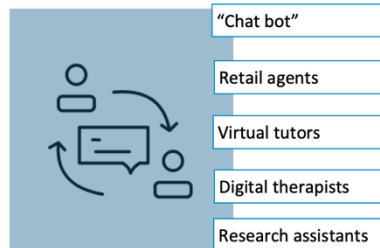------

**System 1: signal detection** *(how the LLM is trained to detect adversarial signals)*

**System 2: signal response** *(how the system triggers actions following detection)*

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

28

# What does this mean for evaluation?

## Conversation



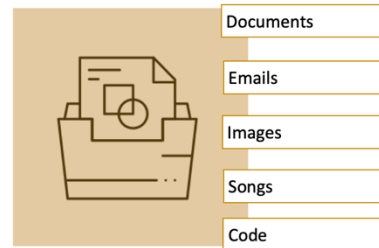| "Chat bot" |
| Retail agents |
| Virtual tutors |
| Digital therapists |
| Research assistants |

**The *thing* is the interaction.**

**Success** = the quality of the discourse.

**Control panel** = pragmatic *fluency* – "co-constructed meaning" with LLM.

**Meaningful evaluation** requires metapragmatic considerations across the discourse frame. *(aka context)*
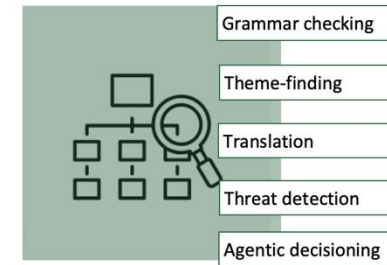
## Generation



| Documents |
| Emails |
| Images |
| Songs |
| Code |

**The *thing* is the artifact.**

**Success** = the quality of the delivery.

**Control panel** = the UI to input criteria. It may or may not involve language.

**Meaningful evaluation** must center success criteria defined at the *artifact* level *(like standard HCI eval.)*

## Analysis



| Grammar checking |
| Theme-finding |
| Translation |
| Threat detection |
| Agentic decisioning |

**The *thing* is the signal...**
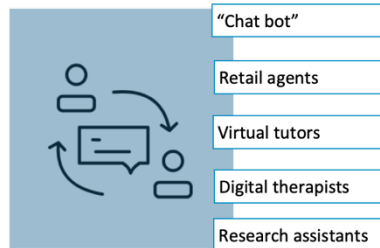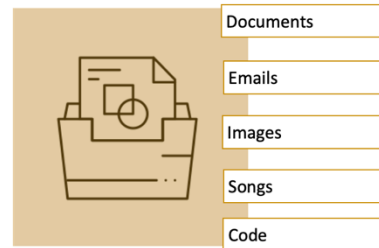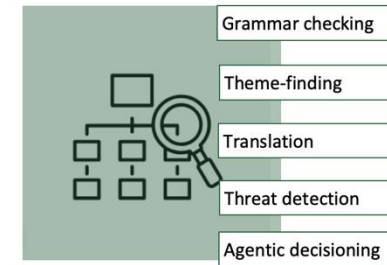**... but really, what you *do* with that signal.**

**Meaningful evaluation** *especially* requires distinguishing *accuracy* from *impacts*.

**System 1 (signal-as-detection):** *Accuracy*, tuned by training data, criteria, thresholds.

**System 2 (signal-as-trigger):** *Appropriateness*, determined by system design decisions.

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

29

# What does this mean for evaluation?

## Conversation



| |
|---|
| "Chat bot" |
| Retail agents |
| Virtual tutors |
| Digital therapists |
| Research assistants |

**The *thing* is the interaction.**

**Success =** the quality of the discourse.

**Control panel =** pragmatic *fluency* – "co-constructed meaning" with LLM.

**Meaningful evaluation** requires metapragmatic considerations across the discourse frame. *(aka context)*

## Generation



| |
|---|
| Documents |
| Emails |
| Images |
| Songs |
| Code |

The *thing* is the artifact.

Success = the quality of the delivery.

Control panel = the UI to input criteria. It may or may not involve language.

Meaningful evaluation must center success criteria defined at the *artifact* level *(like standard HCI eval.)*

## Analysis



| |
|---|
| Grammar checking |
| Theme-finding |
| Translation |
| Threat detection |
| Agentic decisioning |

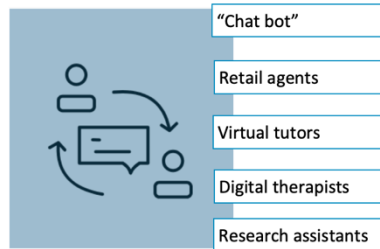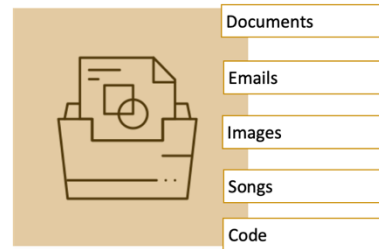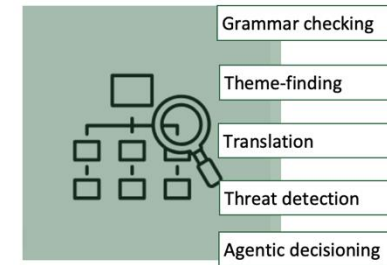The *thing* is the signal...
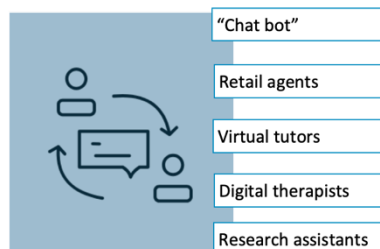... but really, what you *do* with that signal.

Meaningful evaluation *especially* requires distinguishing *accuracy* from *impacts*.

System 1 (signal-as-detection): *Accuracy*, tuned by training data, criteria, thresholds.

System 2 (signal-as-trigger): *Appropriateness*, determined by system design decisions.

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

30

# What does this mean for evaluation?

## Conversation

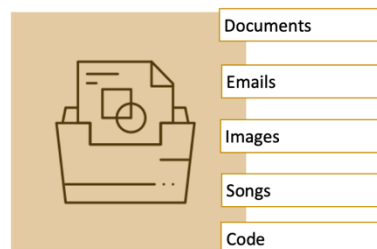| | |
|---|---|
| | "Chat bot" |
| | Retail agents |
| | Virtual tutors |
| | Digital therapists |
| | Research assistants |

The *thing* is the interaction.

**Success** = the quality of the discourse.

**Control panel** = pragmatic *fluency* – "co-constructed meaning" with LLM.

**Meaningful evaluation** requires metapragmatic considerations across the discourse frame. *(aka context)*

## Generation

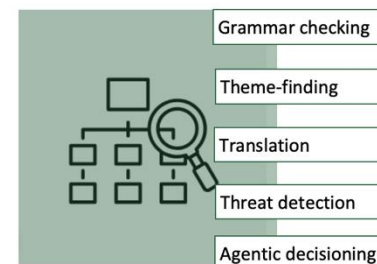| | |
|---|---|
| | Documents |
| | Emails |
| | Images |
| | Songs |
| | Code |

The *thing* is the artifact.

**Success** = the quality of the delivery.

**Control panel** = the UI to input criteria. It may or may not involve language.

**Meaningful evaluation** must center success criteria defined at the *artifact* level *(like standard HCI eval.)*

## Analysis

| | |
|---|---|
| | Grammar checking |
| | Theme-finding |
| | Translation |
| | Threat detection |
| | Agentic decisioning |

The *thing* is the signal…
… but really, what you *do* with that signal.

**Meaningful evaluation** *especially* requires distinguishing *accuracy* from *impacts*.

**System 1 (signal-as-detection):** *Accuracy*, tuned by training data, criteria, thresholds.

**System 2 (signal-as-trigger):** *Appropriateness*, determined by system design decisions.

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

31

# What does this mean for evaluation?

## Conversation



| |
|---|
| "Chat bot" |
| Retail agents |
| Virtual tutors |
| Digital therapists |
| Research assistants |

**The *thing* is the interaction.**

**Success =** the quality of the discourse.

**Control panel =** pragmatic *fluency* – "co-constructed meaning" with LLM.

**Meaningful evaluation** requires metapragmatic considerations across the discourse frame. *(aka context)*

## Generation



| |
|---|
| Documents |
| Emails |
| Images |
| Songs |
| Code |

**The *thing* is the artifact.**

**Success =** the quality of the delivery.

**Control panel =** the UI to input criteria. It may or may not involve language.

**Meaningful evaluation** must center success criteria defined at the *artifact* level *(like standard HCI eval.)*

## Analysis



| |
|---|
| Grammar checking |
| Theme-finding |
| Translation |
| Threat detection |
| Agentic decisioning |

**The *thing* is the signal...**

*... but really, what you do with that signal.*

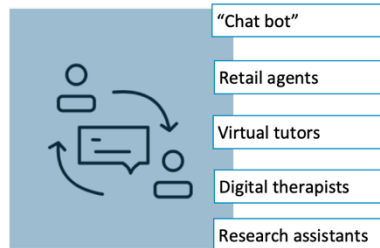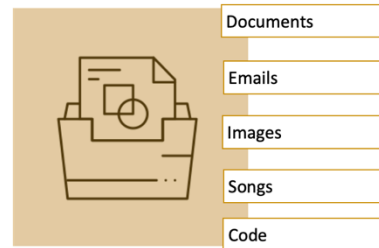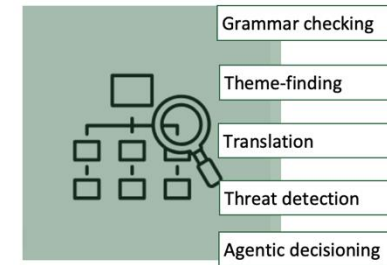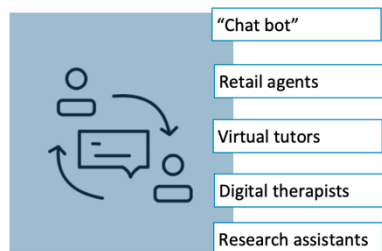*Meaningful evaluation especially requires distinguishing accuracy from impacts.*

*System 1 (signal-as-detection): Accuracy, tuned by training data, criteria, thresholds.*

*System 2 (signal-as-trigger): Appropriateness, determined by system design decisions.*

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

32

# What does this mean for evaluation?

## Conversation



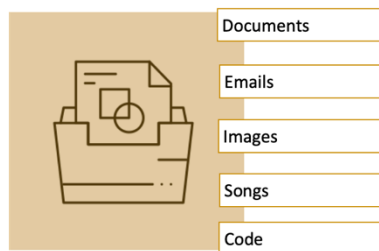| | |
|---|---|
| "Chat bot" | |
| Retail agents | |
| Virtual tutors | |
| Digital therapists | |
| Research assistants | |

**The *thing* is the interaction.**

**Success** = the quality of the discourse.

**Control panel** = pragmatic *fluency* – "co-constructed meaning" with LLM.

**Meaningful evaluation** requires metapragmatic considerations across the discourse frame. *(aka context)*

## Generation



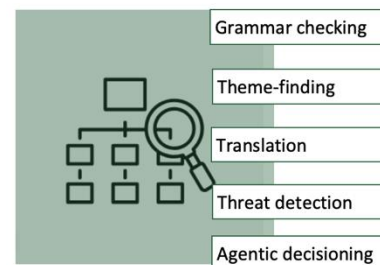| |
|---|
| Documents |
| Emails |
| Images |
| Songs |
| Code |

**The *thing* is the artifact.**

**Success** = the quality of the delivery.

**Control panel** = the UI to input criteria. It may or may not involve language.

**Meaningful evaluation** must center success criteria defined at the *artifact* level *(like standard HCI eval.)*

## Analysis



| |
|---|
| Grammar checking |
| Theme-finding |
| Translation |
| Threat detection |
| Agentic decisioning |

**The *thing* is the signal...**
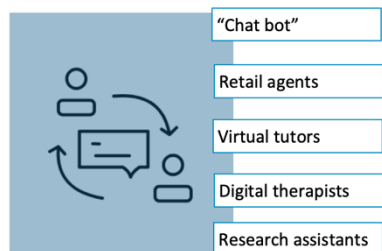**... but really, what you *do* with that signal.**

**Meaningful evaluation** *especially* requires distinguishing *accuracy* from *impacts*.

**System 1 (signal-as-detection):** *Accuracy,* tuned by training data, criteria, thresholds.

**System 2 (signal-as-trigger):** *Appropriateness,* determined by system design decisions.

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

33

# What does this mean for evaluation?

## Conversation



- "Chat bot"
- Retail agents
- Virtual tutors
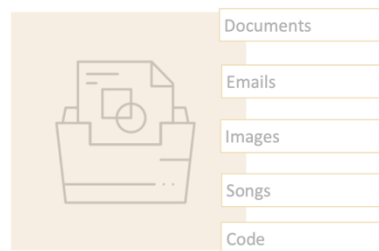- Digital therapists
- Research assistants

**The *thing* is the interaction.**

**Success =** the quality of the discourse.

**Control panel =** pragmatic *fluency* – "co-constructed meaning" with LLM.

**Meaningful evaluation** requires metapragmatic considerations across the discourse frame. *(aka context)*
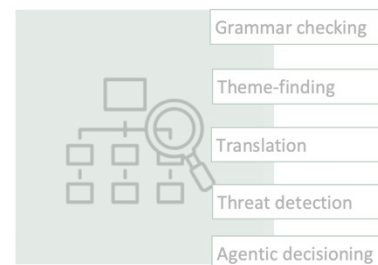
## Generation



- Documents
- Emails
- Images
- Songs
- Code

**The *thing* is the artifact.**

**Success =** the quality of the delivery.

**Control panel =** the UI to input criteria. It may or may not involve language.

**Meaningful evaluation** must center success criteria defined at the *artifact* level *(like standard HCI eval.)*

## Analysis



- Grammar checking
- Theme-finding
- Translation
- Threat detection
- Agentic decisioning

**The *thing* is the signal…**
**… but really, what you *do* with that signal.**

**Meaningful evaluation** *especially* requires distinguishing *accuracy* from *impacts.*

**System 1 (signal-as-detection):** *Accuracy*, tuned by training data, criteria, thresholds.

**System 2 (signal-as-trigger):** *Appropriateness*, determined by system design decisions.

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

34

# Focus: *evaluating conversation*

## Conversation

"Chat bot"

Retail agents

Virtual tutors

Digital therapists

Research assistants

**The *thing* is the interaction.**

**Success =** the quality of the discourse.

**Control panel =** pragmatic *fluency* – "co-constructed meaning" with LLM.

**Meaningful evaluation** requires metapragmatic considerations across the discourse frame. *(aka context)*

## Generation

Documents

Emails

Images

Songs

Code

The *thing* is the artifact.

Success = the quality of the delivery.

Control panel = the UI to input criteria. It may or may not involve language.

Meaningful evaluation must center success criteria defined at the *artifact* level *(like standard HCI eval.)*

## Analysis

Grammar checking

Theme-finding

Translation

Threat detection

Agentic decisioning

The *thing* is the signal...
... but really, what you *do* with that signal.

Meaningful evaluation *especially* requires distinguishing *accuracy* from *impacts*.

System 1 (signal-as-detection): *Accuracy*, tuned by training data, criteria, thresholds.

System 2 (signal-as-trigger): *Appropriateness*, determined by system design decisions.

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

35

# Focus: *evaluating conversation*

Software Engineering Institute

**Evaluating conversational *functionalities* requires applying conversational *methodologies*.**



**DIVIS:**
*Goal: provide environment for victim advocate students to practice leading highly-emotional sexual assault intake interviews*

- First **conversational agent** in the 1960s: ELIZA

- Has since been applied to dozens of different DoD applications

- People apply or adapt their *human-human* language norms to *human-agent* language experiences: **useful for evaluation!**



**PAL3:**
Goal: *on-the-job training*

**Battle buddy:**
Goal: *veteran life quality*

**VITA4Vets:**
Goal: *interviewing skills*

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

36

# Focus: *evaluating conversation*

**Evaluating conversational *functionalities* requires applying conversational *methodologies*.**



**DIVIS:**

*Goal: provide environment for victim advocate students to practice leading highly-emotional sexual assault intake interviews*

- **Discourse analysis:** how to evaluate language in context

  - **Meaning is constructed across multiple turns.**
    *How can you tell if "nice job" sincere praise or sarcastic indictment?*

  - **Roles are explicitly and implicitly negotiated.**
    *Who am I in this conversation? Who are you?*
    *What type of conversation are we having?*

  - **Communication success requires:**
    *Theory of mind: What does this person know?*

    *Grounding & Repair: Given that, what should I say?*
    *Did they know what I mean? How can I get us on the same page?*

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

37

# Focus: *evaluating conversation*

**Evaluating conversational *functionalities* requires applying conversational *methodologies*.**

**This might initially sound like complications...**

**One-shot I/O metrics frequently not realistic**

**Success metrics depend on the *type of discourse***

*Different levels of generative flexibility appropriate for brainstorming vs document summarization.*

**Prioritize "disambiguation" over "accuracy"**

- **Discourse analysis:** how to evaluate language in context

  - **Meaning is constructed across multiple turns.**
    *How can you tell if "nice job" sincere praise or sarcastic indictment?*

  - **Roles are explicitly and implicitly negotiated.**
    *Who am I in this conversation? Who are you?*
    *What type of conversation are we having?*

  - **Communication success requires:**
    *Theory of mind: What does this person know?*

    *Grounding & Repair: Given that, what should I say?*
    *Did they know what I mean? How can I get us on the same page?*
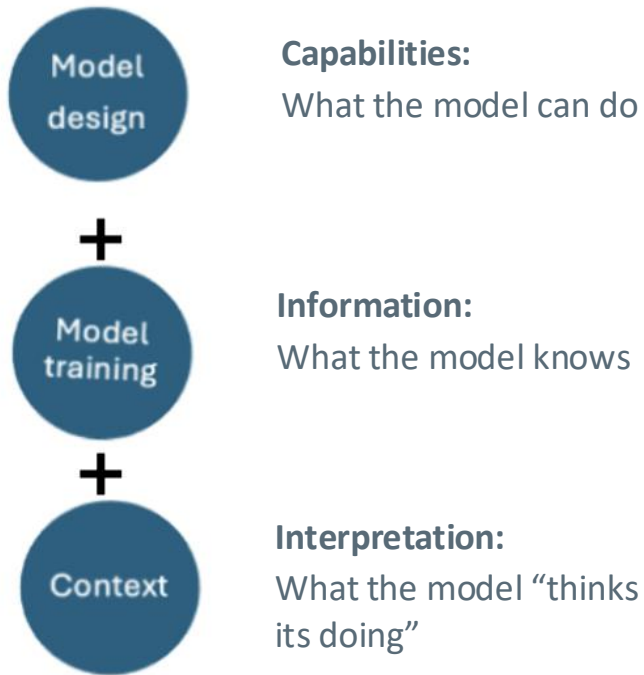
Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

38

# Focus: *evaluating conversation*

**Evaluating conversational *functionalities* requires applying conversational *methodologies*.**

This might initially sound like complications...

One-shot I/O metrics frequently not realistic

Success metrics depend on the *type of discourse*

*Different levels of generative flexibility appropriate for brainstorming vs document summarization.*

Prioritize "disambiguation" over "accuracy"

**...except that we have 60+ years of work to pull from!**

- **Discourse analysis:** how to evaluate language in context

  - **Meaning is constructed across multiple turns.**
    *How can you tell if "nice job" sincere praise or sarcastic indictment?*

  - **Roles are explicitly and implicitly negotiated.**
    *Who am I in this conversation? Who are you?*
    *What type of conversation are we having?*

  - **Communication success requires:**
    *Theory of mind: What does this person know?*

    *Grounding & Repair: Given that, what should I say?*
    *Did they know what I mean? How can I get us on the same page?*

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

Advancing Software for National Security

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

39

# Focus: *evaluating conversation*

**Evaluating conversational *functionalities* requires applying conversational *methodologies*.**

This might initially sound like complications…

One-shot I/O metrics frequently not realistic

Success metrics depend on the *type of discourse*

*Different levels of generative flexibility appropriate for brainstorming vs document summarization.*

Prioritize "disambiguation" over "accuracy"

**…except that we have 60+ years of work to pull from!**

- **Discourse analysis:** how to evaluate language in context

  - **Meaning is constructed across multiple turns.**
    *How can you tell if "nice job" sincere praise or sarcastic indictment?*

  - **Roles are explicitly and implicitly negotiated.**
    *Who am I in this conversation? Who are you?*
    *What type of conversation are we having?*

  - **Communication success requires:**
    *Theory of mind: What does this person know?*

    *Grounding & Repair: Given that, what should I say?*
    *Did they know what I mean? How can I get us on the same page?*

# 1. Discourse is explainable   2. Dialogue is designable

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

40

# Discourse is explainable

*LLM performance requires task alignment across three pillars:*

**Model design**

**+**

**Model training**

**+**

**Context**

**Capabilities:**
What the model can do

**Information:**
What the model knows

**Interpretation:**
What the model "thinks its doing"

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

41

# Discourse is explainable

**LLM performance requires task alignment across three pillars:**



**Capabilities:**
What the model can do

**Information:**
What the model knows

**Interpretation:**
What the model "thinks its doing"

Each of these pillars leads to different *types of problems.*

**Right now, we call them all hallucinations.**

| Term | Definition | Explanation | Mechanism |
|------|-----------|-------------|-----------|
| Interpretive overreach | Premature resolution of ambiguity to maintain conversational flow | Model chooses one plausible interpretation of am ambiguous prompt without surfacing uncertainty or requesting grounding -making its "best guess" given its "understanding" of its context. | Training pressure to maintain local coherence; training emphasis on fluency over epistemic caution. User prompts for being clear and concise may exacerbate this risk. |
| Fictive cohesion | Inserting discursive 'connective tissue' to maintain immediate-turn smoothness at the expense of precision | The model inserts plausible-sounding connections to make responses feel complete. This is often harmless (or desirable) when users are internally consistent and accurate, but otherwise can reify misperceptions | Next-token prediction objective factors well-formed continuity; generative content becomes introduced in the gaps between distinct ideas - sometimes as desirable generative insight, but other times misleading or inaccurate |
| Template overfit | Over-influence of learned genre, role, script, or format structures that conflict with users' distinct task directives | The model defaults to familiar structural patterns, missing user intentions. This is functionally similar to human bias or getting 'too comfortable' in a conversation / falling into old patterns. | Dominance of certain discursive frames in training data; insufficient task disambiguation cues in user prompt |
| Dramaturgical loyalty | Output aligned with LLMs "perceived" communicative role, incentivizing undesirable performance goals | The model tailors responses to match and inferred social script, such as mirroring tone or rhetorical strategy - (this leads to desirable cooperative alignment when effective) | Context-sensitive fine-tuning, human feedback that emphasizes politeness, helpfulness, or friendliness norms, especially in default (non-interventionist) prompt spaces. |
| Source collapse | Misattribution or source blending | The LLM "loses track" of where specific information comes from, merging voices or claims | Distributional encoding of semantically similar text can lead to blending due to the probabilistic nature of LLM outputs |
| Epistemic misalignment | Failure to mark an ideated or hypothesized connection with appropriate stance markers | The problem isn't the generative output, it's the failure to successfully convey the nature of that output. (e.g., a metaphor that isn't understood to be a metaphor is a lie). | Pressure for concision can override stance-marking conventions; Human-LLM misunderstanding about the meaning of certain stance cues; LLM 'losing track' of the stance |

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

42

# Dialogue is designable

| LLM performance quality on any given task is shaped by the level of alignment across three distinct pillars | | | Without these distinctions, we are less able to target strategic impact pathways. | |
|---|---|---|---|---|
| Pillar: | ...matters, because: | ...directly impacts... | When failure is driven by inappropriate... | ... the relevant impact pathways involve: |
| **Model design** | The model's:<br>• Tokenization scheme<br>• Context window<br>• Positional encoding | • Semantic distortions<br>• Multi-turn coherence<br>• Order sensitivity | capability<br><br>(misalignment of task x model design) | • Building new models<br>• Task breakdown<br>• Tool supplementation |
| **Model training** | The model's:<br>• Training objective<br>• Data coverage<br>• Fine-tuning strategy | • Statistical associations<br>• Representations<br>• Default interpretations | knowledge<br><br>(misalignment of task x model training) | • Collecting more data<br>• Retraining models<br>• RAG / memory augmentation |
| **Task cueing** | The user's:<br>• Lexical indexing<br>• Framing / task criteria<br>• Metapragmatics | • Interpretive indexing<br>• Task prioritization<br>• Dramaturgical role | interpretation<br><br>(misalignment of task x cueing) | • Turn structuring<br>• Discourse visibility<br>• Grounding initiations |

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

43

# Dialogue is designable

| LLM performance quality on any given task is shaped by the level of alignment across three distinct pillars | | | Without these distinctions, we are less able to target strategic impact pathways. | |
|---|---|---|---|---|
| **Pillar:** | **...matters, because:** | **...directly impacts...** | When failure is driven by inappropriate... | ... the relevant impact pathways involve: |
| **Model design** | The model's:<br>• Tokenization scheme<br>• Context window<br>• Positional encoding | • Semantic distortions<br>• Multi-turn coherence<br>• Order sensitivity | capability<br>(misalignment of task x model design) | • Building new models<br>• Task breakdown<br>• Tool supplementation |
| **Model training** | The model's:<br>• Training objective<br>• Data coverage<br>• Fine-tuning strategy | • Statistical associations<br>• Representations<br>• Default interpretations | knowledge<br>(misalignment of task x model training) | • Collecting more data<br>• Retraining models<br>• RAG / memory augmentation |
| **Task cueing** | The user's:<br>• Lexical indexing<br>• Framing / task criteria<br>• Metapragmatics | • Interpretive indexing<br>• Task prioritization<br>• Dramaturgical role | interpretation<br>(misalignment of task x cueing) | • Turn structuring<br>• Discourse visibility<br>• Grounding initiations |

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

44

# Dialogue is designable

| Pillar: | ...matters, because: | ...directly impacts... | When failure is driven by inappropriate... | ... the relevant impact pathways involve: |
|---|---|---|---|---|
| **Model design** | The model's:<br>• Tokenization scheme<br>• Context window<br>• Positional encoding | • Semantic distortions<br>• Multi-turn coherence<br>• Order sensitivity | capability<br>(misalignment of task x model design) | • Building new models<br>• Task breakdown<br>• Tool supplementation |
| **Model training** | The model's:<br>• Training objective<br>• Data coverage<br>• Fine-tuning strategy | • Statistical associations<br>• Representations<br>• Default interpretations | knowledge<br>(misalignment of task x model training) | • Collecting more data<br>• Retraining models<br>• RAG / memory augmentation |
| **Task cueing** | The user's:<br>• Lexical indexing<br>• Framing / task criteria<br>• Metapragmatics<br>**Prompt?** | • Interpretive indexing<br>• Task prioritization<br>• Dra... ...ble<br>**Output?** | interpretation<br>(misalignment of task x cueing) | • Turn structuring<br>• Discourse visibility<br>• Grounding initiations |

*LLM performance quality on any given task is shaped by the level of alignment across three distinct pillars*

*Without these distinctions, we are less able to target strategic impact pathways.*

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

45

# Dialogue is designable

| | LLM performance quality on any given task is shaped by the level of alignment across three distinct pillars | | Without these distinctions, we are less able to target strategic impact pathways. | |
|---|---|---|---|---|
| Pillar: | ...matters, because: | ...directly impacts... | When failure is driven by inappropriate... | ... the relevant impact pathways involve: |
| **Model design** | The model's:<br>• Tokenization scheme<br>• Context window<br>• Positional encoding | • Semantic distortions<br>• Multi-turn coherence<br>• Order sensitivity | **capability**<br>(misalignment of task x model design) | • Building new models<br>• Task breakdown<br>• Tool supplementation |
| **Model training** | The model's:<br>• Training objective<br>• Data coverage<br>• Fine-tuning strategy | • Statistical associations<br>• Representations<br>• Default interpretations | **knowledge**<br>(misalignment of task x model training) | • Collecting more data<br>• Retraining models<br>• RAG / memory augmentation |
| **Task cueing** | The user's:<br>• Lexical indexing<br>• Framing / task criteria<br>• Metapragmatics | • Interpretive indexing<br>• Task prioritization<br>• Dramaturgical role | **interpretation**<br>(misalignment of task x cueing) | • Turn structuring<br>• Discourse visibility<br>• Grounding initiations |

**Prompt?**

**Output?**

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

46

# Dialogue is designable

| LLM performance quality on any given task is shaped by the level of alignment across three distinct pillars | | | Without these distinctions, we are less able to target strategic impact pathways. | |
|---|---|---|---|---|
| **Pillar:** | **...matters, because:** | **...directly impacts...** | **When failure is driven by inappropriate...** | **... the relevant impact pathways involve:** |
| **Model design** | The model's:<br>• Tokenization scheme<br>• Context window<br>• Positional encoding | • Semantic distortions<br>• Multi-turn coherence<br>• Order sensitivity | **capability**<br>(misalignment of task x model design) | • Building new models<br>• Task breakdown<br>• Tool supplementation |
| **Model training** | The model's:<br>• Training objective<br>• Data coverage<br>• Fine-tuning strategy | • Statistical associations<br>• Representations<br>• Default interpretations | **knowledge**<br>(misalignment of task x model training) | • Collecting more data<br>• Retraining models<br>• RAG / memory augmentation |
| **Task cueing** | The user's:<br>• Lexical indexing<br>• Framing / task criteria<br>• Metapragmatics | • Interpretive indexing<br>• Task prioritization<br>• Dramaturgical role | **interpretation**<br>(misalignment of task x cueing) | • Turn structuring<br>• Discourse visibility<br>• Grounding initiations |

**Prompt?**

**Output?**

**Prompt engineering?**

**Prompt better?**

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

47

# Dialogue is designable

| Pillar: | ...matters, because: | ...directly impacts... | When failure is driven by inappropriate... | ... the relevant impact pathways involve: |
|---|---|---|---|---|
| | LLM performance quality on any given task is shaped by the level of alignment across three distinct pillars | | Without these distinctions, we are less able to target strategic impact pathways. | |
| **Model design** | The model's: <br>• Tokenization scheme <br>• Context window <br>• Positional encoding | • Semantic distortions <br>• Multi-turn coherence <br>• Order sensitivity | **capability** <br><br>(misalignment of task x model design) | • Building new models <br>• Task breakdown <br>• Tool supplementation |
| **Model training** | The model's: <br>• Training objective <br>• Data coverage <br>• Fine-tuning strategy | • Statistical associations <br>• Representations <br>• Default interpretations | **knowledge** <br><br>(misalignment of task x model training) | • Collecting more data <br>• Retraining models <br>• RAG / memory augmentation |
| **Task cueing** | The user's: <br>• Lexical indexing <br>• Framing / task criteria <br>• Metapragmatics | • Interpretive indexing <br>• Task prioritization <br>• Dramaturgical role | interpretation <br>(misalignment of task x cueing) <br>**Prompt engineering?** | • Turn structuring <br>• Discourse visibility <br>• Grounding initiations <br>**Prompt better?** |

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

48

# Dialogue is designable

| | LLM performance quality on any given task is shaped by the level of alignment across three distinct pillars | | Without these distinctions, we are less able to target strategic impact pathways. | |
|---|---|---|---|---|
| Pillar: | ...matters, because: | ...directly impacts... | When failure is driven by inappropriate... | ... the relevant impact pathways involve: |
| **Model design** | The model's:<br>• Tokenization scheme<br>• Context window<br>• Positional encoding | • Semantic distortions<br>• Multi-turn coherence<br>• Order sensitivity | **capability**<br>(misalignment of task x model design) | • Building new models<br>• Task breakdown<br>• Tool supplementation |
| **Model training** | The model's:<br>• Training objective<br>• Data coverage<br>• Fine-tuning strategy | • Statistical associations<br>• Representations<br>• Default interpretations | **knowledge**<br>(misalignment of task x model training) | • Collecting more data<br>• Retraining models<br>• RAG / memory augmentation |
| **Task cueing** | The user's:<br>• Lexical indexing<br>• Framing / task criteria<br>• Metapragmatics | • Interpretive indexing<br>• Task prioritization<br>• Dramaturgical role | **interpretation**<br>(misalignment of task x cueing) | • Turn structuring<br>• Discourse visibility<br>• Grounding initiations |

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

49

# Discourse analysis as a diagnostic lens

## 1. **Identify** what type of discourse failure occurred

- *Surface where dialogue breaks down (e.g., interpretive overreach, failed implicature…)*
- *Treat chat logs as structured evidence – rich dataset that you have for free!*

## 2. **Explain** the mechanism behind the misalignment

- *Surface potential hypotheses that explain patterns from established literature*
- *Experiments: contrastive trials contrasting sociolinguistic explanations*
- *Center explanatory mechanism in the science of evaluation*

## 3. **Intervene** at the level of the interaction

- *Design and test discourse-level fixes (scaffolds, stance markers, role visibility).*
- *Evaluate not just system accuracy, but alignment and coordination gains.*

Title of the Presentation Goes Here
© 2025 Carnegie Mellon University

**Advancing Software for National Security**

[DISTRIBUTION STATEMENT Please copy and paste the appropriate distribution statement into this space.]

50