# Natural Language Processing

1 December 2025

# Briefly, in the time remaining:

Some project information

Course wrap-up and reflection

# Project information

| | | | |
|---|---|---|---|
| Mon. | Dec. | 1 | You are here |
| Wed. | Dec. | 3 | **Project presentations** |
| | | | |
| Mon. | Dec. | 8 | **Project presentations** |
| Wed. | Dec. | 10 | Last day of classes – but not for us; it's a "Friday" |
| | | | 11:59 p.m.   **Project draft due** |
| | | | |
| Sun. | Dec. | 14 | End of study period |
| | | | 11:59 p.m.   **Project due** |

# Presentations

Describing work in progress

Think of it as a chance to tell me and your classmates the cool things you're working on – and to get some feedback and questions while you still have time to address them!

Team of 1–2:  5–6 minutes
Team of 3:     Up to 8 minutes

We want to have maximum time for presentations and minimize the time switching between groups, so I'll have all the slides ready to go on one computer.

Please send me your slides (Google Slides, PDF, PowerPoint, Keynote) before the start of the class where you're presenting.

# What have we learned?

"The idea of giving computers the ability to process human language is as old as the idea of computers themselves."

Jurafsky & Martin, 2nd ed.

Tasks

Techniques

| Tasks | | Techniques |
|-------|---|-----------|

*Tokenization*

| Tasks | Techniques |
|-------|------------|
| Tokenization | Regular expressions |

| Tasks | Techniques |
|-------|------------|
| *Tokenization* | *Regular expressions* |
| | *Rule-based systems* |

| *Tasks* | *Techniques* | |
|---|---|---|
| Tokenization | Regular expressions | BPE |
| | Rule-based systems | |

| *Tasks* | *Techniques* |
| --- | --- |
| *Tokenization* | *Regular expressions*     *BPE* |
| *Sentence segmentation* | *Rule-based systems* |

| Tasks | Techniques |
|---|---|
| | |

Tasks

Tokenization

Sentence segmentation

Representation learning

Techniques

Regular expressions          BPE

Rule-based systems

| Tasks | Techniques |
|-------|-----------|
| Tokenization | Regular expressions     BPE |
| Sentence segmentation | Rule-based systems |
| Representation learning | Word count embeddings |

| Tasks | | Techniques | |
|---|---|---|---|
| *Tokenization* | | *Regular expressions* | *BPE* |
| *Sentence segmentation* | | *Rule-based systems* | |
| *Representation learning* | | *Word count embeddings* | *Word2vec* |

| Tasks | Techniques |
|-------|-----------|

**Tasks**

Tokenization

Sentence segmentation

Representation learning

Classification
*Sentiment analysis, author identification, etc.*

**Techniques**

Regular expressions          BPE

Rule-based systems

Word count embeddings        Word2vec

| Tasks | Techniques |
|---|---|
| Tokenization | Regular expressions      BPE |
| Sentence segmentation | Rule-based systems |
| Representation learning | Word count embeddings      Word2vec |
| Classification<br>*Sentiment analysis, author identification, etc.* | Logistic regression |

| Tasks | Techniques | |
|---|---|---|
| *Tokenization* | *Regular expressions* | *BPE* |
| *Sentence segmentation* | *Rule-based systems* | |
| *Representation learning* | *Word count embeddings* | *Word2vec* |
| *Classification*<br>*Sentiment analysis, author identification, etc.* | *Logistic regression* | *Neural networks* |

| Tasks | Techniques |
|---|---|

**Tasks**

Tokenization

Sentence segmentation

Representation learning

Classification
*Sentiment analysis, author identification, etc.*

Language modeling

**Techniques**

Regular expressions          BPE

Rule-based systems

Word count embeddings          Word2vec

Logistic regression          Neural networks

| Tasks | Techniques |
|-------|-----------|
| | |

Tasks

Techniques

Tokenization

Regular expressions          BPE

Sentence segmentation

Rule-based systems

Representation learning

Word count embeddings     Word2vec

Classification
*Sentiment analysis, author identification, etc.*

Logistic regression      Neural networks

Language modeling

N-grams

| Tasks | Techniques |
|---|---|
| | |

**Tasks**

Tokenization

Sentence segmentation

Representation learning

Classification
*Sentiment analysis, author identification, etc.*

Language modeling

**Techniques**

Regular expressions          BPE

Rule-based systems

Word count embeddings          Word2vec

Logistic regression          Neural networks

N-grams          Transformers

| Tasks | Techniques | |
|---|---|---|
| Tokenization | Regular expressions | BPE |
| Sentence segmentation | Rule-based systems | |
| Representation learning | Word count embeddings | Word2vec |
| Classification<br>*Sentiment analysis, author identification, etc.* | Logistic regression | Neural networks |
| Language modeling | N-grams | Transformers |
| Text generation | | |

| *Tasks* | *Techniques* | |
|---|---|---|
| Tokenization | Regular expressions | BPE |
| Sentence segmentation | Rule-based systems | |
| Representation learning | Word count embeddings | Word2vec |
| Classification<br>*Sentiment analysis, author identification, etc.* | Logistic regression | Neural networks |
| Language modeling | N-grams | Transformers |
| Text generation | Chaining | |

| Tasks | Techniques | |
|---|---|---|
| Tokenization | Regular expressions | BPE |
| Sentence segmentation | Rule-based systems | |
| Representation learning | Word count embeddings | Word2vec |
| Classification _Sentiment analysis, author identification, etc._ | Logistic regression | Neural networks |
| Language modeling | N-grams | Transformers |
| Text generation | Chaining | Autoregressive models |

# Skills

*Text processing*: Working with strings, reading in text datasets, tokenizing, cleaning text, regular expressions

*Evaluation*: Computing accuracy, precision, recall, F-score, perplexity, and cosine similarity
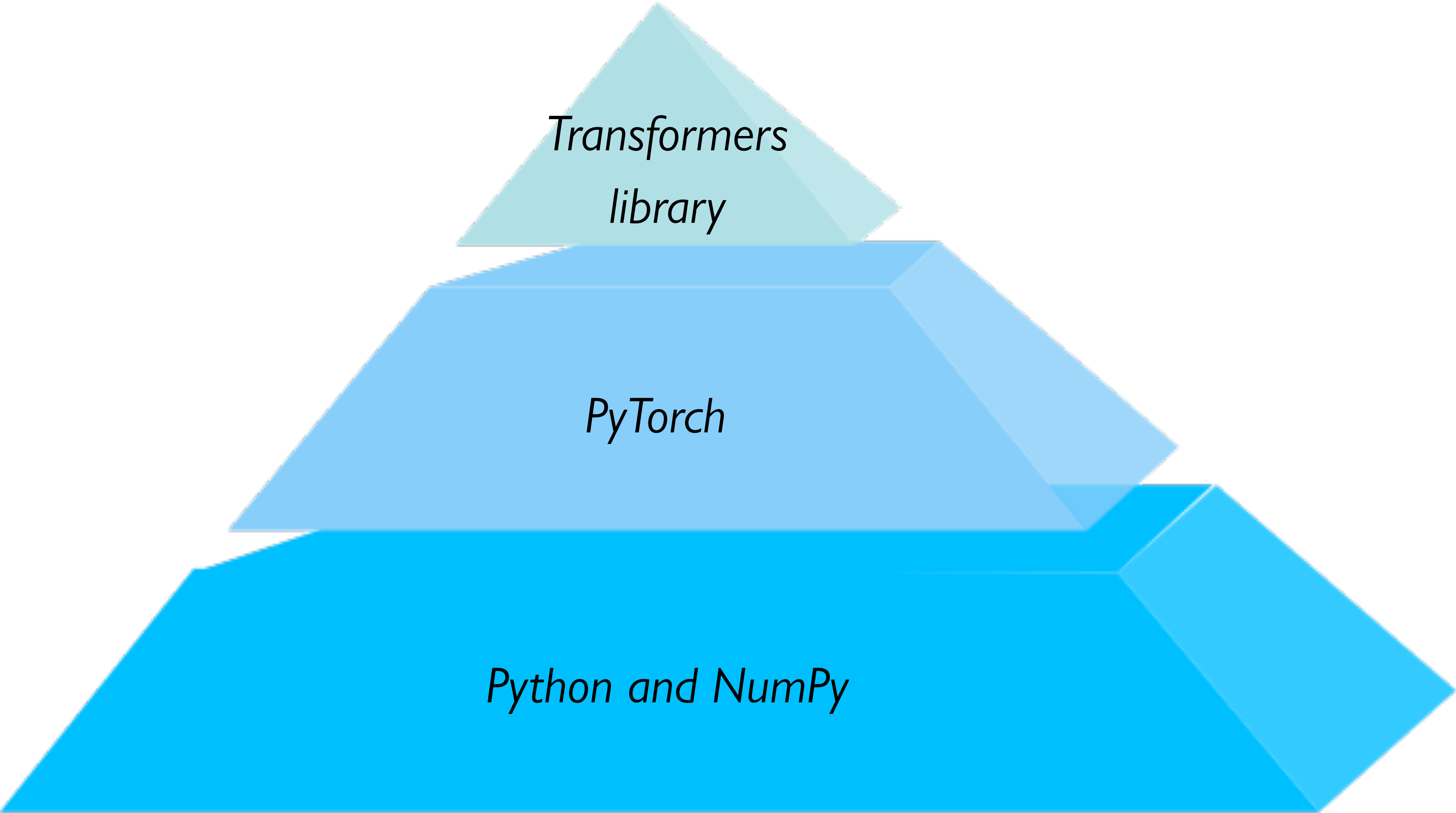
*Implementing classification models*

*Building* and *running* neural network models
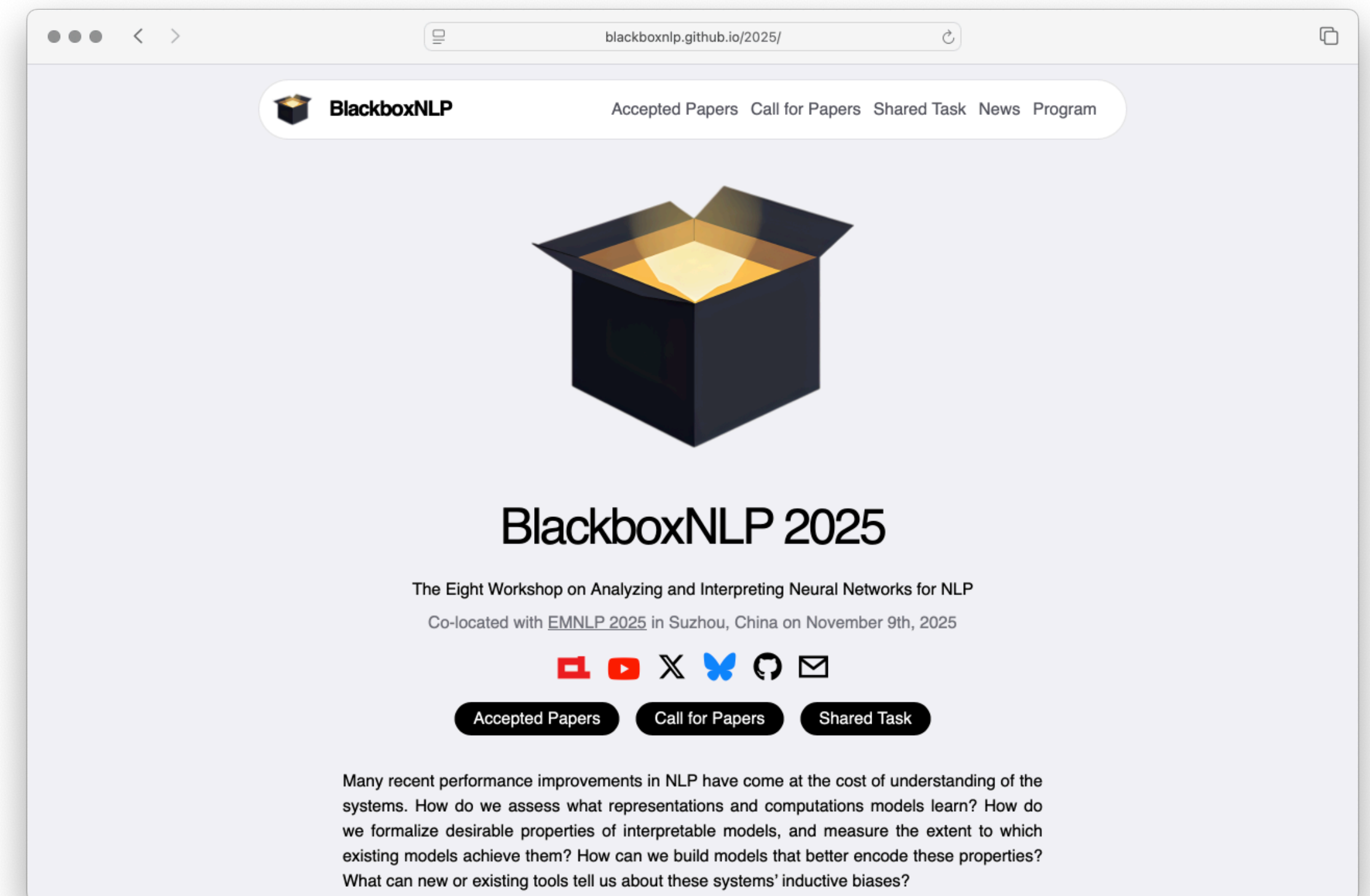
More abstraction
& collaboration

Transformers
library

PyTorch

Python and NumPy

More control &
understanding

What is there still to learn?

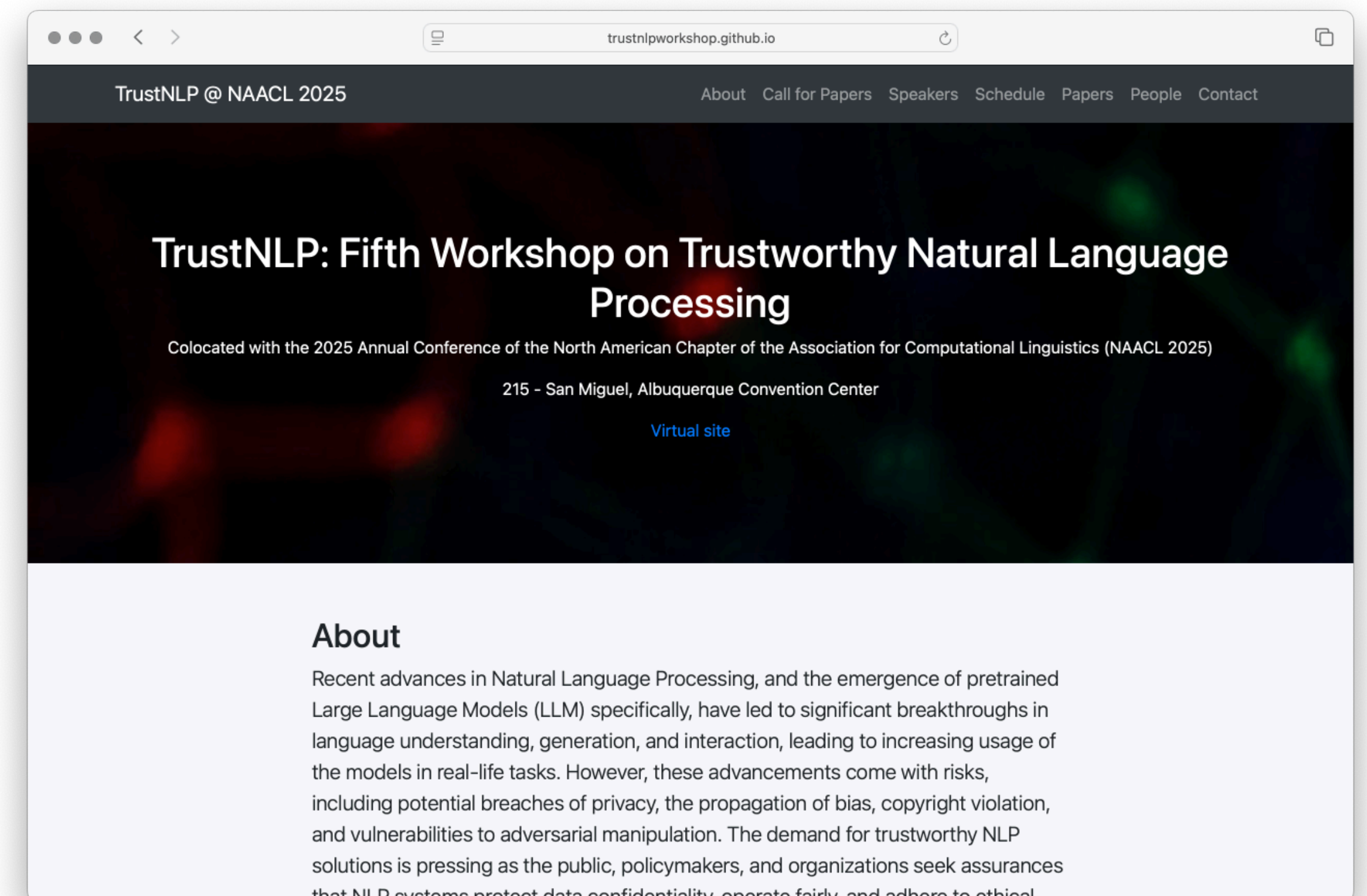# Interpretability and explainability

How do blackbox models make decisions?

What are large language models learning, and how do they use this information?

What are the limits of what
large language models can do?
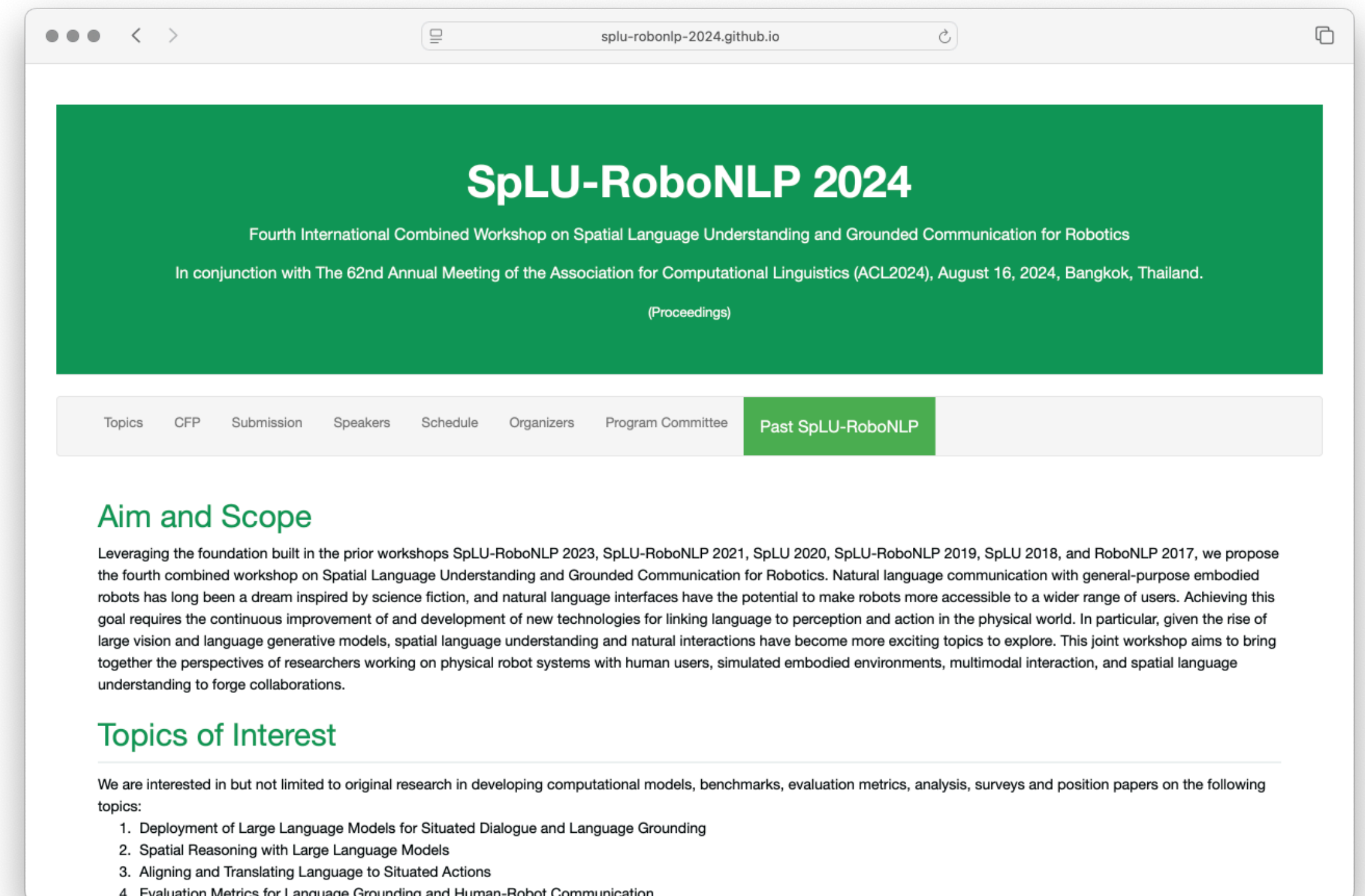
How can we intervene when
models are incorrect?
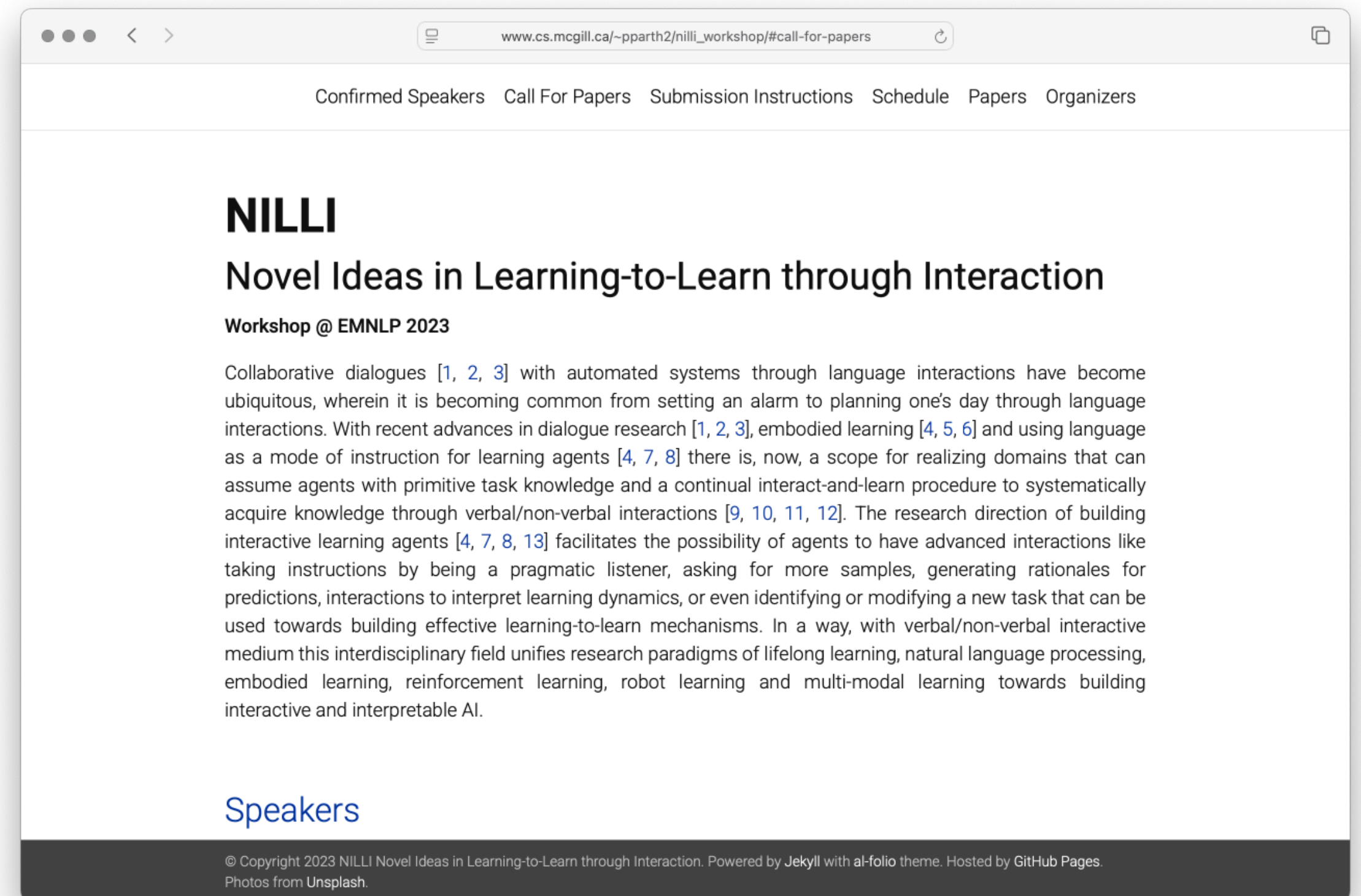
# Grounding and multimodality

How do we combine text-based models with other modalities, like vision and perceptual information?

How can we help models relate to the physical world?

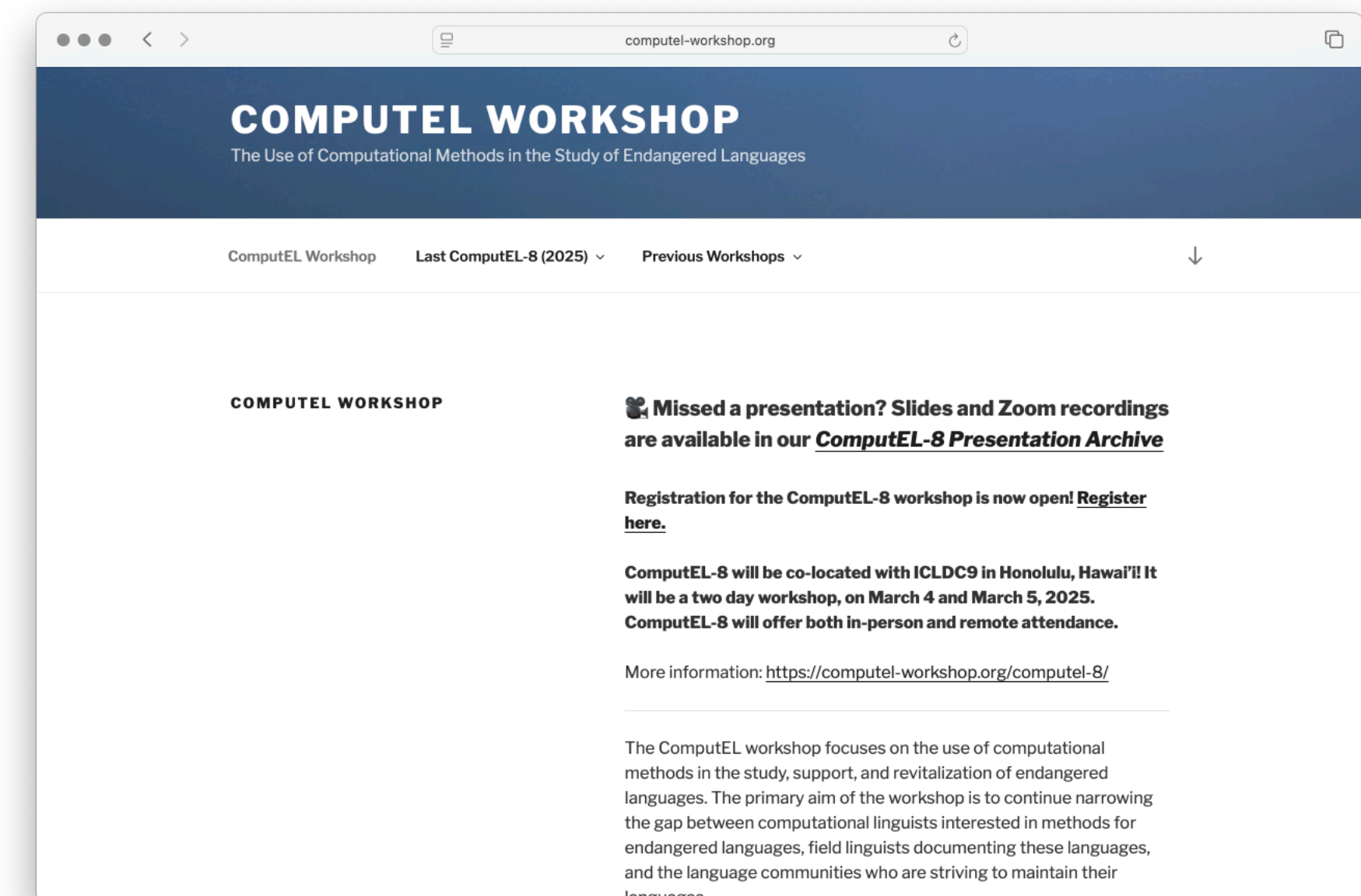How do we help models understand social interaction?

How much of a role does access to perceptual information play in learning?

Confirmed Speakers   Call For Papers   Submission Instructions   Schedule   Papers   Organizers

# NILLI
## Novel Ideas in Learning-to-Learn through Interaction
**Workshop @ EMNLP 2023**

Collaborative dialogues [1, 2, 3] with automated systems through language interactions have become ubiquitous, wherein it is becoming common from setting an alarm to planning one's day through language interactions. With recent advances in dialogue research [1, 2, 3], embodied learning [4, 5, 6] and using language as a mode of instruction for learning agents [4, 7, 8] there is, now, a scope for realizing domains that can assume agents with primitive task knowledge and a continual interact-and-learn procedure to systematically acquire knowledge through verbal/non-verbal interactions [9, 10, 11, 12]. The research direction of building interactive learning agents [4, 7, 8, 13] facilitates the possibility of agents to have advanced interactions like taking instructions by being a pragmatic listener, asking for more samples, generating rationales for predictions, interactions to interpret learning dynamics, or even identifying or modifying a new task that can be used towards building effective learning-to-learn mechanisms. In a way, with verbal/non-verbal interactive medium this interdisciplinary field unifies research paradigms of lifelong learning, natural language processing, embodied learning, reinforcement learning, robot learning and multi-modal learning towards building interactive and interpretable AI.
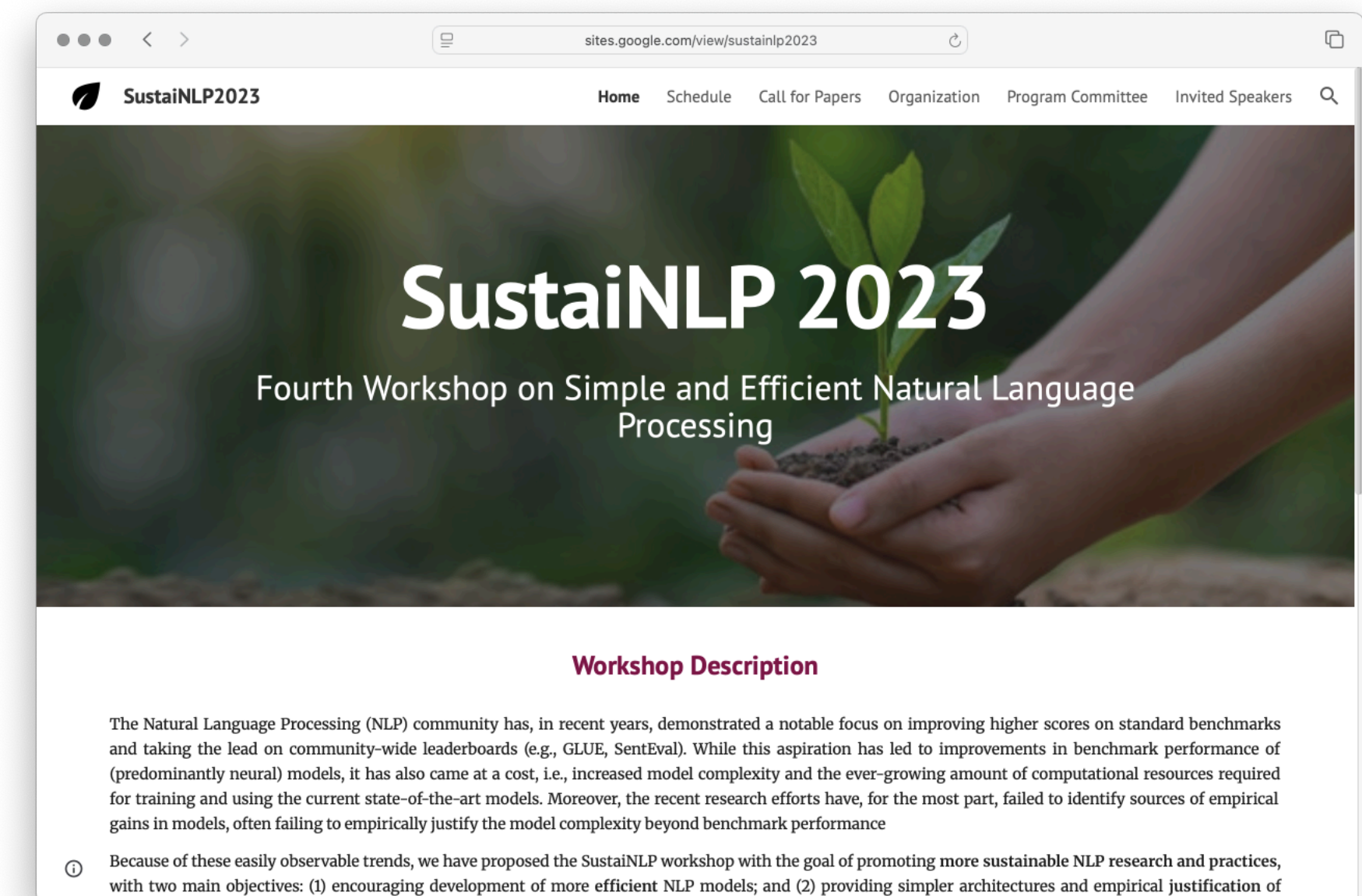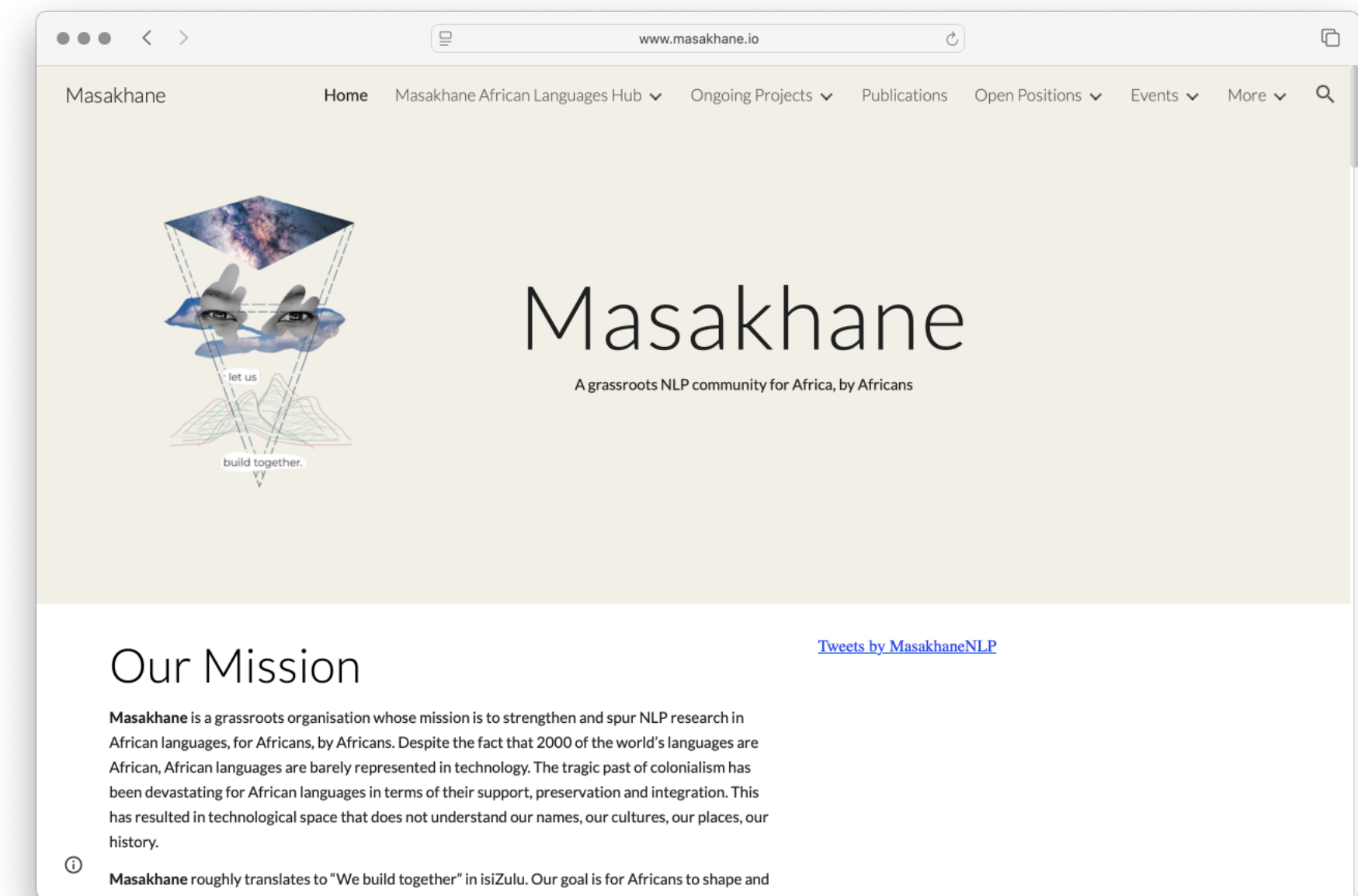
## Speakers

# Linguistic diversity

How can we improve NLP tools for low-resource languages?

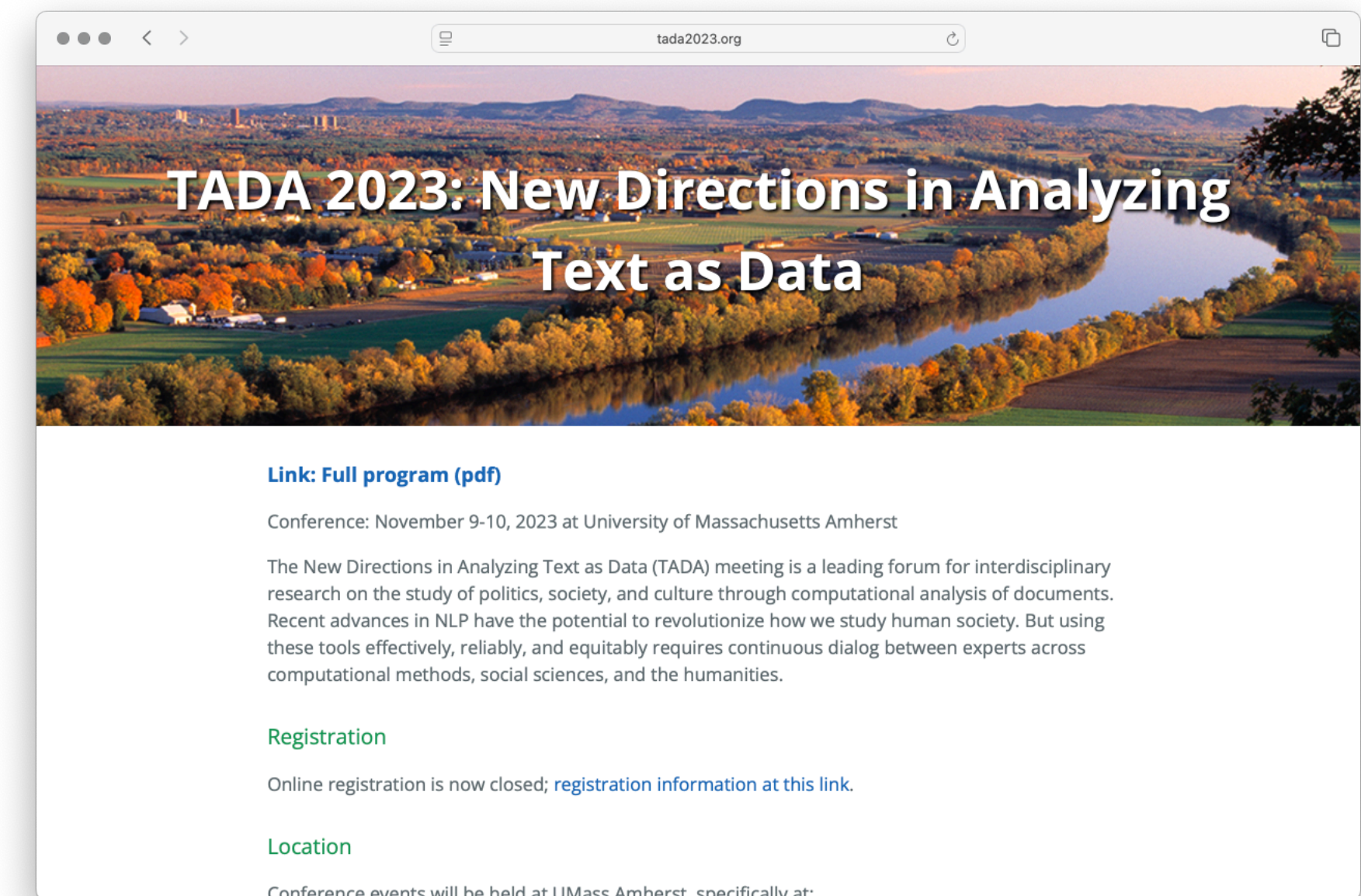How do models deal with different language varieties?

How can we achieve good performance with sparse data? With small models and limited compute?

# Computational social science

How can we share our best approaches for text processing with other fields?

How well do NLP techniques work with limited data? In niche domains? In tasks that require expert judgment?

go.vassar.edu/course/evals

# Acknowledgments

The class incorporates material from:

Carolyn Anderson, Wellesley College