# Prompting Is Programming: A Query Language for Large Language Models

**Author: Luca et al.**
**Presentor: Dongjae Lee**

**20240514**

Luca Beurer-Kellner, Marc Fischer, & Martin T. Vechev (2023). Prompting Is Programming: A Query Language for Large Language Models. Proc. ACM Program. Lang., 7(PLDI), 1946–1969.
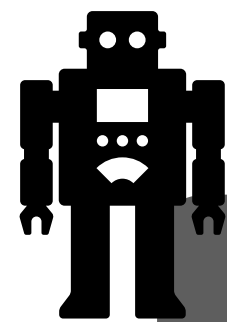
# Problem of interaction with language model

- Human in the loop!

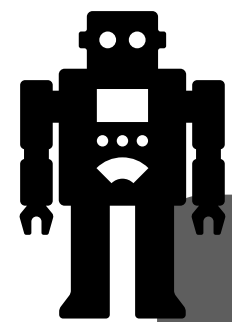  - Grunt work

# Problem of interaction with language model

- Human in the loop!

  - Grunt work


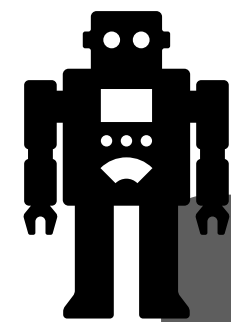
1 + 1 = 3

# Problem of interaction with language model

- Human in the loop!

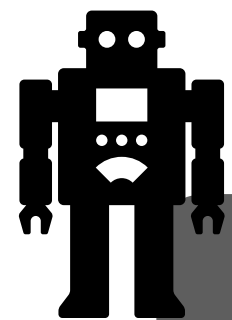  - Grunt work

1 + 1 = 3

😡

# Problem of interaction with language model

- Human in the loop!

  - Grunt work

1 + 1 = 3

😠

1 + 1 = 2

# Problem of interaction with language model

- Human in the loop!

  - Grunt work

# Problem of interaction with language model
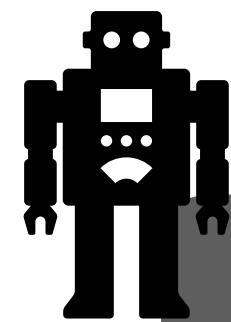
- Human in the loop!

  - Grunt work

1 + 1 = 3

😡

1 + 1 = 2

😁

**Require too much labor**

# Problem of interaction with language model

- Algorithmic way

  - Require lots of reimplementation!

# Problem of interaction with language model

- Algorithmic way

  - Require lots of reimplementation!



1 + 1 = 3

# Problem of interaction with language model

- Algorithmic way

  - Require lots of reimplementation!



1 + 1 = 3

😡

# Problem of interaction with language model

- Algorithmic way

  - Require lots of reimplementation!



1 + 1 = 3

😡

1 + 1 = 2

# Problem of interaction with language model

- Algorithmic way

  - Require lots of reimplementation!

1 + 1 = 3

😠

1 + 1 = 2

😁

# Problem of interaction with language model

- Algorithmic way

  - Require lots of reimplementation!

def fibonacci(n): ~~~

???

😡

😢

# Problem of interaction with language model

- Algorithmic way

  - Require lots of reimplementation!



def fibonacci(n): ~~~

???

😡

😢

**Constraint checker must be reimplemented from the scratch**

# What is the cause of the problem?

- No abstraction and interface!

# What is the cause of the problem?

- No abstraction and interface!

- Computer Science has very well-designed abstraction layers

# What is the cause of the problem?

- No abstraction and interface!

- Computer Science has very well-designed abstraction layers

  - Hardware (Architecture) / Software

# What is the cause of the problem?

- No abstraction and interface!

- Computer Science has very well-designed abstraction layers

  - Hardware (Architecture) / Software

  - Assembly / Programming Language

# What is the cause of the problem?

- No abstraction and interface!

- Computer Science has very well-designed abstraction layers

  - Hardware (Architecture) / Software

  - Assembly / Programming Language

  - ...

# What is the cause of the problem?

- No abstraction and interface!

- Computer Science has very well-designed abstraction layers

  - Hardware (Architecture) / Software

  - Assembly / Programming Language

  - ...

- Language Models

# What is the cause of the problem?

- No abstraction and interface!

- Computer Science has very well-designed abstraction layers

  - Hardware (Architecture) / Software

  - Assembly / Programming Language

  - ...

- Language Models

  - No **interface** between language models and applications

# What is the cause of the problem?

- No abstraction and interface!

- Computer Science has very well-designed abstraction layers

  - Hardware (Architecture) / Software

  - Assembly / Programming Language

  - ...

- Language Models

  - No **interface** between language models and applications

  - **Then... Let's make it!**

# LMQL

# LMQL

What is your question?

My question is ...

Answer: ...

😁

**LMQL Script**

# LMQL



What is your question?

My question is …

Answer: …

😄

**LMQL Script**

**LMQL interpreter**

# LMQL

What is your question?

My question is ...

Answer: ...

😁

LMQL Script

LMQL interpreter

Language Model

# Interface for language model

- **Langauge model programming (LMP)**

# Interface for language model

- **Langauge model programming (LMP)**

  - Abstract language model prompting and interaction

# Interface for language model

- **Langauge model programming (LMP)**

  - Abstract language model prompting and interaction

  - Combining text prompting and scripting (some algorithm)

# Interface for language model

- **Langauge model programming (LMP)**

  - Abstract language model prompting and interaction

  - Combining text prompting and scripting (some algorithm)

- **Language Model Query Language (LMQL)**

# Interface for language model

- **Langauge model programming (LMP)**

  - Abstract language model prompting and interaction

  - Combining text prompting and scripting (some algorithm)

- **Language Model Query Language (LMQL)**

  - **Abstraction of language model**

# Interface for language model

- **Langauge model programming (LMP)**

  - Abstract language model prompting and interaction

  - Combining text prompting and scripting (some algorithm)

- **Language Model Query Language (LMQL)**

  - **Abstraction of language model**

  - We can write down **scripts for interacting** with a language model

# Interface for language model

- **Langauge model programming (LMP)**

  - Abstract language model prompting and interaction

  - Combining text prompting and scripting (some algorithm)

- **Language Model Query Language (LMQL)**

  - **Abstraction of language model**

  - We can write down **scripts for interacting** with a language model

  - We can give **some constraints** at the high level (word, number, ...)

# How to use LMQL

- LMQL is Python + SQL!

- Outline of LMQL Program

# How to use LMQL

- LMQL is Python + SQL!

- Outline of LMQL Program

```
<decoder> <query>

from <model>

[where <cond>]

[distribute <dist>]
```

# How to use LMQL

- LMQL is Python + SQL!

- Outline of LMQL Program

```
<decoder> <query>

from <model>

[where <cond>]

[distribute <dist>]
```

Python code

# Example

**LMQL Program**

```
argmax

   a = 2

   b = 3

   "What is {a} + {b}? [c]"

   "The answer is {c}"
from "gpt2-large"
where INT(c)
```

# Example

**LMQL Program**

```
argmax                    We are going to use greedy decoding

   a = 2

   b = 3

   "What is {a} + {b}? [c]"

   "The answer is {c}"
from "gpt2-large"
where INT(c)
```

# Example

**LMQL Program**

```
argmax

  a = 2

  b = 3

  "What is {a} + {b}? [c]"

  "The answer is {c}"

from "gpt2-large"

where INT(c)
```

Use language model gpt2-large

# Example

**LMQL Program**

```
argmax

    a = 2

    b = 3

    "What is {a} + {b}? [c]"

    "The answer is {c}"

from "gpt2-large"

where INT(c)
```

Variable c must be integer

# Example

**LMQL Program**

**Environment**

**Prompt**

""

```
argmax

  a = 2

  b = 3

  "What is {a} + {b}? [c]"

  "The answer is {c}"

from "gpt2-large"

where INT(c)
```

# Example

**LMQL Program**              **Environment**          **Prompt**

■■■■■■■                    `""`

`argmax`

    `a = 2` ——————→ Add variable **a** to environment

    `b = 3`

    `"What is {a} + {b}? [c]"`

    `"The answer is {c}"`

`from "gpt2-large"`

`where INT(c)`

# Example

### LMQL Program

**Environment**

**Prompt**

a = 2

""

```
argmax

  a = 2

  b = 3

  "What is {a} + {b}? [c]"

  "The answer is {c}"

from "gpt2-large"

where INT(c)
```

# Example

**LMQL Program**                                    **Environment**              **Prompt**

a = 2                                          ""

```
argmax

   a = 2

   b = 3  ──────────────▶  Add variable b to environment

   "What is {a} + {b}? [c]"

   "The answer is {c}"

from "gpt2-large"

where INT(c)
```

# Example

### LMQL Program

**argmax**

   a = 2

   b = 3

   "What is {a} + {b}? [c]"

   "The answer is {c}"

**from** "gpt2-large"

**where** INT(c)

### Environment

```
a = 2
b = 3
```

### Prompt

""

# Example

**LMQL Program**　　　　　　　　　　　**Environment**　　　　**Prompt**

```
a = 2
b = 3
```

""

```
argmax

    a = 2

    b = 3

    "What is {a} + {b}? [c]"  ⟶

    "The answer is {c}"

from "gpt2-large"

where INT(c)
```

**{?}**: Get **?** from environment

**[c]**: Get response from LM.
　　　Then, save it in environment

# Example

**LMQL Program**

**Environment**

**Prompt**

| a = 2 | "What is 2 + 3? " |
| b = 3 | |

```
argmax

    a = 2

    b = 3

    "What is {a} + {b}? [c]"

    "The answer is {c}"

from "gpt2-large"

where INT(c)
```

# Example

**LMQL Program**

**Environment**

**Prompt**

```
a = 2
b = 3
```

"What is 2 + 3? "

**argmax**

```
  a = 2

  b = 3

  "What is {a} + {b}? [c]"

  "The answer is {c}"
from "gpt2-large"

where INT(c)
```

Prompt: "**What is 2 + 3? **"

gpt2-large: "**5**"

17

# Example

### LMQL Program

**Environment**  **Prompt**

```
argmax

    a = 2

    b = 3

    "What is {a} + {b}? [c]"

    "The answer is {c}"

from "gpt2-large"

where INT(c)
```

a = 2
b = 3
c = 5

"What is 2 + 3? 5"

# Example

**LMQL Program**

**Environment**

**Prompt**

```
a = 2
b = 3
c = 5
```

"What is 2 + 3? 5"

**argmax**

   a = 2

   b = 3

   "What is {a} + {b}? [c]"

   "The answer is {c}" ⟶ **{?}**: Get **?** from environment

**from** "gpt2-large"

**where** INT(c)

# Example

### LMQL Program

**Environment**

**Prompt**

```
a = 2
b = 3
c = 5
```

"What is 2 + 3? 5
The answer is 5"

```
argmax

    a = 2

    b = 3

    "What is {a} + {b}? [c]"

    "The answer is {c}"

from "gpt2-large"

where INT(c)
```

# Key features

- Interaction with **control flow and environment**

# Key features

- Interaction with **control flow and environment**

  - We can use **branch, loop, and function calls** like programming language

# Key features

- Interaction with **control flow and environment**

  - We can use **branch, loop, and function calls** like programming language

  - We can use **variable**

# Key features

- Interaction with **control flow and environment**

  - We can use **branch, loop, and function calls** like programming language

  - We can use **variable**

- Constraint decoding

# Key features

- Interaction with **control flow and environment**

  - We can use **branch, loop, and function calls** like programming language

  - We can use **variable**

- Constraint decoding

  - We can give some constraints through `where` statement

# Key features

- Interaction with **control flow and environment**

  - We can use **branch, loop, and function calls** like programming language

  - We can use **variable**

- Constraint decoding

  - We can give some constraints through `where` statement

  - You can **specify all constraints** that can be expressed with Python Expr

# Key features

- Interaction with **control flow and environment**

  - We can use **branch, loop, and function calls** like programming language

  - We can use **variable**

- <mark>Constraint decoding</mark>

  - We can give some constraints through `where` statement

  - You can **specify all constraints** that can be expressed with Python Expr

# What is constraint decoding?

- Generate output to satisfy the given constraints

# What is constraint decoding?

- Generate output to satisfy the given constraints

- Constraints can be

# What is constraint decoding?

- Generate output to satisfy the given constraints

- Constraints can be

  - Format (Syntax)

# What is constraint decoding?

- Generate output to satisfy the given constraints

- Constraints can be

  - Format (Syntax)

  - Type

# What is constraint decoding?

- Generate output to satisfy the given constraints

- Constraints can be

  - Format (Syntax)

  - Type

  - Length of output

# What is constraint decoding?

- Generate output to satisfy the given constraints

- Constraints can be

  - Format (Syntax)

  - Type

  - Length of output

  - Termination Condition

# What is constraint decoding?

- Generate output to satisfy the given constraints

- Constraints can be

  - Format (Syntax)

  - Type

  - Length of output

  - Termination Condition

  - …

# Naïve approach

- Check **after** decode

# Naïve approach

- Check **after** decode

# Naïve approach

- Check **after** decode



**Language
Model**

# Naïve approach

- Check **after** decode



Language
Model

# Naïve approach

- Check **after** decode



Language
Model

Constraint
Checker

# Naïve approach

- Check **after** decode



Language
Model

Constraint
Checker

# Naïve approach

- Check **after** decode



Language
Model

Constraint
Checker

# Naïve approach

- Check **after** decode



Language Model

Constraint Checker

# Naïve approach

- Check **after** decode



Language Model

Constraint Checker

# Naïve approach

- Check **after** decode



- Problem: **Decoding time >>> Checking time**

# Efficient approach

- Check **before** decode

# Efficient approach

- Check **before** decode

# Efficient approach

- Check **before** decode



**Constraint Checker**

# Efficient approach

- Check **before** decode



Constraint
Checker

Which token can
come next?

# Efficient approach

- Check **before** decode



Constraint
Checker

Feasible
Token Set

Which token can
come next?

# Efficient approach

- Check **before** decode

# Efficient approach

- Check **before** decode

# Efficient approach

- Check **before** decode

# Efficient approach

- Check **before** decode



- Feasible Token Set is always finite?

# Efficient approach

- Check **before** decode



Which token can come next?

**Constraint Checker**

**Feasible Token Set**

**Language Model**

- Feasible Token Set is always finite?

  - Yes

# Efficient approach

- Check **before** decode



- Feasible Token Set is always finite?

    - Yes

    - Language Models have a finite token set

# How to make a feasible token set?

- 2 Types of token

# How to make a feasible token set?

- 2 Types of token

  - **Maybe** satisfying the constraint

# How to make a feasible token set?

- 2 Types of token

  - **Maybe** satisfying the constraint

  - **Never** satisfying the constraint

# How to make a feasible token set?

- 2 Types of token

  - **Maybe** satisfying the constraint

  - **Never** satisfying the constraint

- We cannot 100% guarantee that some tokens satisfy the given constraint

# How to make a feasible token set?

- 2 Types of token

  - **Maybe** satisfying the constraint

  - **Never** satisfying the constraint

- We cannot 100% guarantee that some tokens satisfy the given constraint

- But we can guarantee that some tokens **never satisfy** the constraint

# 2 Types of token

# 2 Types of token

Constraint: `eval(Var) < 10 and match(Var, "[0-9+]*")`

# 2 Types of token

Constraint: `eval(Var) < 10 and match(Var, "[0-9+]*")`

Current Output: **2+0+1+4+??**

# 2 Types of token

Constraint: **eval(Var) < 10 and match(Var, "[0-9+]\*")**

Current Output: **2+0+1+4+??**

Maybe satisfying: **0, 1, 2**

# 2 Types of token

Constraint: `eval(Var) < 10 and match(Var, "[0-9+]*")`

Current Output: `2+0+1+4+??`

Maybe satisfying: `0, 1, 2` $\longrightarrow$ Depending on the trailing tokens

# 2 Types of token

Constraint: **eval**(**Var**) < **10** **and match**(**Var**, **"[0-9+]*"**)

Current Output: **2+0+1+4+??**

Maybe satisfying: **0, 1, 2** ⟶ Depending on the trailing tokens

Never satisfying: **3, 4, ...**

# 2 Types of token

Constraint: `eval(Var) < 10 and match(Var, "[0-9+]*")`

Current Output: `2+0+1+4+??`

Maybe satisfying: `0, 1, 2` $\longrightarrow$ Depending on the trailing tokens

Never satisfying: `3, 4, ...`

Focus on removing **"never satisfying"** tokens

# Partial evalution

- To check the constraint, we should evaluate the current **unfinished output**

# Partial evalution

- To check the constraint, we should evaluate the current **unfinished output**

- **2 Semantics** for partial evaluation

# Partial evalution

- To check the constraint, we should evaluate the current **unfinished output**

- **2 Semantics** for partial evaluation

  - Final Semantics

# Partial evalution

- To check the constraint, we should evaluate the current **unfinished output**

- **2 Semantics** for partial evaluation

  - Final Semantics

    - Denote current expression can be changed or not

# Partial evalution

- To check the constraint, we should evaluate the current **unfinished output**

- **2 Semantics** for partial evaluation

  - Final Semantics

    - Denote current expression can be changed or not

    - 4 States: FIN, VAR, INC, DEC

# Partial evalution

- To check the constraint, we should evaluate the current **unfinished output**

- **2 Semantics** for partial evaluation

  - Final Semantics

    - Denote current expression can be changed or not

    - 4 States: FIN, VAR, INC, DEC

  - Follow Map

# Partial evalution

- To check the constraint, we should evaluate the current **unfinished output**

- **2 Semantics** for partial evaluation

  - Final Semantics

    - Denote current expression can be changed or not

    - 4 States: FIN, VAR, INC, DEC

  - Follow Map

    - Denote how the current state change according to the trailing token

# Final Semantics

- Evaluate expression as 4 states

# Final Semantics

- Evaluate expression as 4 states

  - FIN: The evaluation result of expression never changing

# Final Semantics

- Evaluate expression as 4 states

  - FIN: The evaluation result of expression never changing

  - VAR: The evaluation result of expression changing

# Final Semantics

- Evaluate expression as 4 states

  - FIN: The evaluation result of expression never changing

  - VAR: The evaluation result of expression changing

  - INC: The evaluation result of expression keeps increasing

# Final Semantics

- Evaluate expression as 4 states

  - FIN: The evaluation result of expression never changing

  - VAR: The evaluation result of expression changing

  - INC: The evaluation result of expression keeps increasing

    - len(output)

# Final Semantics

- Evaluate expression as 4 states

  - FIN: The evaluation result of expression never changing

  - VAR: The evaluation result of expression changing

  - INC: The evaluation result of expression keeps increasing

    - len(output)

  - DEC: The evaluation result of expression keeps decreasing

# Final Semantics

| expression | $\text{FINAL}[\,\cdot\,;\sigma]$ |
|---|---|
| $\langle const \rangle$ | FIN |
| python variable $\langle pyvar \rangle$ | VAR |
| previous hole $\langle var \rangle$ | FIN |
| current var $\langle var \rangle$ | INC |
| future hole $\langle var \rangle$ | INC |
| words($v$) | $\text{FINAL}[v]$ |
| sentences($v$) | $\text{FINAL}[v]$ |
| len($v$) | $\text{FINAL}[v]$ |

number equality $n == m$
$$\begin{cases} \text{FIN} & \text{if } \text{FINAL}[n] = \text{FIN} \\ & \wedge \text{FINAL}[m] = \text{FIN} \\ \text{VAR} & \text{else} \end{cases}$$

string equality $x == y$
$$\begin{cases} \text{FIN} & \text{if } \text{FINAL}[x] = \text{FIN} \\ & \wedge \text{FINAL}[y] = \text{FIN} \\ \text{FIN} & \exists i \bullet x[i] \neq y[i] \\ & \wedge \text{FINAL}[x] \neq \text{VAR} \\ & \wedge \text{FINAL}[y] \neq \text{VAR} \\ \text{VAR} & \text{else} \end{cases}$$

function $\text{fn}(\tau_1, \ldots, \tau_k)$
$$\begin{cases} \text{FIN} & \text{if } \bigwedge_{i=1}^{k} a(\tau_i) = \text{FIN} \\ \text{VAR} & \text{else} \end{cases}$$

| expression | $\text{FINAL}[\,\cdot\,;\sigma]$ |
|---|---|

stop_at(var, $s$)
$$\begin{cases} \text{FIN} & \text{if } [\![var]\!]_\sigma.\text{endswith}(s) \\ & \wedge \text{FINAL}[var] = \text{INC} \\ \text{VAR} & \text{else} \end{cases}$$

$x$ in $s$ for strings $x, s$
$$\begin{cases} \text{FIN} & \text{if } x \text{ in } s \wedge \text{FINAL}[x] = \text{FIN} \\ & \wedge \text{FINAL}[s] = \text{INC} \\ \text{VAR} & \text{else} \end{cases}$$

$e$ in $l$ for string $e$, set $l$
$$\begin{cases} \text{FIN} & \text{if } \nexists i \in l \bullet i.\text{startswith}(e) \\ & \wedge \text{FINAL}[x] \in \{\text{INC, FIN}\} \\ & \wedge \text{FINAL}[l] = \text{FIN} \\ \text{VAR} & \text{else} \end{cases}$$

$x < y$
$$\begin{cases} \text{FIN} & \text{if } x{<}y \wedge \text{FINAL}[x] \in \{\text{DEC, FIN}\} \\ & \wedge \text{FINAL}[y] \in \{\text{INC, FIN}\} \\ \text{VAR} & \text{else} \end{cases}$$

$a$ and $b$
$$\begin{cases} \text{FIN} & \text{if } \exists v \in \{a, b\} \bullet [\![v]\!]_\sigma^F = \text{FIN}(\bot) \\ \text{FIN} & \text{if } \forall v \in \{a, b\} \bullet [\![v]\!]_\sigma^F = \text{FIN}(\top) \\ \text{VAR} & \text{else} \end{cases}$$

$a$ or $b$
$$\begin{cases} \text{FIN} & \text{if } \exists v \in \{a, b\} \bullet [\![v]\!]_\sigma^F = \text{FIN}(\top) \\ \text{FIN} & \text{if } \forall v \in \{a, b\} \bullet [\![v]\!]_\sigma^F = \text{FIN}(\bot) \\ \text{VAR} & \text{else} \end{cases}$$

| not $a$ | $\text{FINAL}[a]$ |
|---|---|

# Follow Map

- Follow Map denotes **state transition** according to the trailing token

# Follow Map

- Follow Map denotes **state transition** according to the trailing token

  - State: Final semantics's state + Possible values

# Follow Map

- Follow Map denotes **state transition** according to the trailing token

  - State: Final semantics's state + Possible values

$$\text{FOLLOW}\big[\text{TEXT in ["Stephen Hawking"]}\big](\text{"Steph"}, t) = \begin{cases} \text{FIN}(\top) & \text{if } t = \text{"en Hawking"} \\ \text{FIN}(\bot) & \text{else} \end{cases}$$

# Follow Map

- Follow Map denotes **state transition** according to the trailing token

  - State: Final semantics's state + Possible values

$$\text{Follow}\big[\text{TEXT in } [\text{"Stephen Hawking"}]\big](\text{"Steph"}, t) = \begin{cases} \text{FIN}(\top) & \text{if } t = \text{"en Hawking"} \\ \text{FIN}(\bot) & \text{else} \end{cases}$$

Constraint

# Follow Map

- Follow Map denotes **state transition** according to the trailing token

  - State: Final semantics's state + Possible values

$$\text{Follow}\big[\textsf{TEXT in } [\texttt{"Stephen Hawking"}]\big](\texttt{"Steph"}, t) = \begin{cases} \text{FIN}(\top) & \text{if } t = \texttt{"en Hawking"} \\ \text{FIN}(\bot) & \text{else} \end{cases}$$

Constraint    Current Output

# Follow Map

- Follow Map denotes **state transition** according to the trailing token

  - State: Final semantics's state + Possible values

$$\text{FOLLOW}\big[\text{TEXT in }[\texttt{"Stephen Hawking"}]\big](\texttt{"Steph"}, t) = \begin{cases} \text{FIN}(\top) & \text{if } t = \texttt{"en Hawking"} \\ \text{FIN}(\bot) & \text{else} \end{cases}$$

Constraint     Current Output     Transitions

# Follow Map

- Follow Map denotes **state transition** according to the trailing token

  - State: Final semantics's state + Possible values

$$\text{FOLLOW}\big[\text{TEXT in } [\text{"Stephen Hawking"}]\big](\text{"Steph"}, t) = \begin{cases} \text{FIN}(\top) & \text{if } t = \text{"en Hawking"} \\ \text{FIN}(\bot) & \text{else} \end{cases}$$

Constraint      Current Output      Transitions

- If the trailing token is "en Hawking", the evaluation result is **True (Top)** and **never changes (FIN)**

30

# Follow Map

- Follow Map denotes **state transition** according to the trailing token

  - State: Final semantics's state + Possible values

$$\text{Follow}\big[\text{TEXT in ["Stephen Hawking"]}\big](\text{"Steph"}, t) = \begin{cases} \text{FIN}(\top) & \text{if } t = \text{"en Hawking"} \\ \text{FIN}(\bot) & \text{else} \end{cases}$$

Constraint     Current Output     Transitions

- If the trailing token is "en Hawking", the evaluation result is **True (Top)** and **never changes (FIN)**

- Else, the evaluation result is **False (Bottom)** and **never changes (FIN)**

# Follow Map

| expression | FOLLOW$[\cdot](u, t)$ |
|---|---|
| $\langle const \rangle$ | $[\![\langle const \rangle]\!]_\sigma$ |
| python variable $\langle pyvar \rangle$ | $[\![pyvar]\!]_{\sigma[v \leftarrow vt]}$ |
| previous hole $\langle var \rangle$ | $[\![\langle var \rangle]\!]_\sigma$ |
| current var $v$ | $\begin{cases} \text{FIN}(v) & \text{if } t = \text{EOS} \\ \text{INC}(vt) & \text{else} \end{cases}$ |
| future hole $\langle var \rangle$ | None |
| words$(v)$ | $\begin{cases} \text{FIN}(w_1, \ldots, w_k) & \text{if } t = \text{EOS} \\ \text{INC}(w_1, \ldots, w_k) & \text{if } t = {}_\sqcup \\ \text{INC}(w_1, \ldots, w_k\, t) & \text{else} \end{cases}$ where $w_1, \ldots, w_k \leftarrow [\![words(v)]\!]_\sigma$ |
| sentences$(v)$ | $\begin{cases} \text{FIN}(s_1, \ldots, s_k) & \text{if } t = \text{EOS} \\ \text{INC}(s_1, \ldots, s_k, t) & \text{if } s_k.\text{endswith(".")} \\ \text{INC}(s_1, \ldots, s_k\, t) & \text{else} \end{cases}$ where $s_1, \ldots, s_k \leftarrow [\![sentences(v)]\!]_\sigma$ |
| len$(v)$ | $\begin{cases} \text{len}(v) & \text{if } t = \text{EOS} \\ \text{len}(v) + 1 & \text{else} \end{cases}$ |
| len$(l)$ over list $l$ | $len([\![l]\!]_{\sigma[v \leftarrow vt]})$ |

| expression | FOLLOW$[\cdot](u, t)$ |
|---|---|
| fn$(\tau_1, \ldots, \tau_k)$ | fn$([\![\tau_1]\!]_{\sigma[v \leftarrow vt]}, \ldots, [\![\tau_k]\!]_{\sigma[v \leftarrow vt]})$ |
| stop_at$(var, s)$ | $\begin{cases} \text{FIN}(b) & \text{if } b \wedge \text{FINAL}[var] = \text{INC} \\ \text{VAR}(l) & \text{else} \end{cases}$ where $b = [\![var]\!]_\sigma.\text{endswith}(s)$ |
| x in $s$ for string $s$ and constant $x$ | $\begin{cases} \top & \text{if x in s} \vee \text{x in } t \\ \bot & \text{else} \end{cases}$ |
| x in $l$ for constant list/set $l$ | $\begin{cases} \text{FIN}(\top) & \text{if t in l} \\ \text{VAR}(\bot) & \text{if } \exists e \in l \bullet \\ & \quad e.\text{startswith}(vt) \\ \bot & \text{else} \end{cases}$ |
| x < y | $[\![x]\!]_{\sigma[v \leftarrow vt]} < [\![y]\!]_{\sigma[v \leftarrow vt]}$ |
| string comp. a == $v$ | $\begin{cases} \text{FIN}(\top) & \text{if } vt = a \\ \text{VAR}(\bot) & \text{if a.startswith}(vt) \\ \bot & \text{else} \end{cases}$ |
| number comp. x == y | $[\![x]\!]_{\sigma[v \leftarrow vt]} = [\![y]\!]_{\sigma[v \leftarrow vt]}$ |
| a and b | $[\![x]\!]_{\sigma[v \leftarrow vt]}$ and $[\![y]\!]_{\sigma[v \leftarrow vt]}$ |
| a or b | $[\![x]\!]_{\sigma[v \leftarrow vt]}$ or $[\![y]\!]_{\sigma[v \leftarrow vt]}$ |
| not a | not $[\![x]\!]_{\sigma[v \leftarrow vt]}$ |

# Soundness

- How to define soundness?

# Soundness

- How to define soundness?

- Soundness: **Removed tokens are must <span style="color:red">unsafe</span>**

# Soundness

- How to define soundness?

- Soundness: **Removed tokens are must <span style="color:red">unsafe</span>**

$$\forall t \in \mathcal{V} \bullet (\text{FOLLOW}[e])(u, t) = \text{FIN}(\bot) \Rightarrow [\![e]\!]_{\sigma[v \leftarrow ut]} = \text{FIN}(\bot)$$

# Soundness

- How to define soundness?

- Soundness: **Removed tokens are must <span style="color:red">unsafe</span>**

$$\forall t \in \mathcal{V} \bullet (\textsc{Follow}[e])(u, t) = \textsc{Fin}(\bot) \Rightarrow [\![e]\!]_{\sigma[v \leftarrow ut]} = \textsc{Fin}(\bot)$$

If the result of the state according to the follow map semantics is **false for all tokens**

# Soundness

- How to define soundness?

- Soundness: **Removed tokens are must <span style="color:red">unsafe</span>**

$$\forall t \in \mathcal{V} \bullet (\text{Follow}[e])(u, t) = \text{FIN}(\bot) \Rightarrow [\![e]\!]_{\sigma[v \leftarrow ut]} = \text{FIN}(\bot)$$

**Concrete evaluation result must be false**

# Soundness

- How to define soundness?

- Soundness: **Removed tokens are must <span style="color:red">unsafe</span>**

$$\forall t \in \mathcal{V} \bullet (\text{Follow}[e])(u, t) = \text{FIN}(\bot) \Rightarrow [\![e]\!]_{\sigma[v \leftarrow ut]} = \text{FIN}(\bot)$$

**Concrete evaluation result must be false**

- It means the left tokens are still unsafe

# Soundness

- How to define soundness?

- Soundness: **Removed tokens are must <span style="color:red">unsafe</span>**

$$\forall t \in \mathcal{V} \bullet (\text{FOLLOW}[e])(u, t) = \text{FIN}(\bot) \Rightarrow \llbracket e \rrbracket_{\sigma[v \leftarrow ut]} = \text{FIN}(\bot)$$

**Concrete evaluation result must be false**

- It means the left tokens are still unsafe

    - Backtracking!

33

# Evaluation

- 3 Research Question

# Evaluation

- 3 Research Question

  - Expressiveness: Can we implement advanced prompting techniques?

# Evaluation

- 3 Research Question

  - Expressiveness: Can we implement advanced prompting techniques?

  - Performance: Can LMQL lower the number of model queries?

# Evaluation

- 3 Research Question

    - Expressiveness: Can we implement advanced prompting techniques?

    - Performance: Can LMQL lower the number of model queries?

    - Accuracy: Does LMQL's constraint decoding affect task accuracy?

# Evaluation

- 3 Research Question

  - Expressiveness: Can we implement advanced prompting techniques?

  - Performance: Can LMQL lower the number of model queries?

  - Accuracy: Does LMQL's constraint decoding affect task accuracy?

- Baseline

# Evaluation

- 3 Research Question

  - Expressiveness: Can we implement advanced prompting techniques?

  - Performance: Can LMQL lower the number of model queries?

  - Accuracy: Does LMQL's constraint decoding affect task accuracy?

- Baseline

  - Huggingface's generate() API

# Expressiveness

- Chain-of-Thought Prompting

```
argmax
    "Pick the odd word out: skirt, dress, pen, jacket."
    "skirt is clothing, dress is clothing, pen is an object, jacket is clothing."
    "So the odd one is pen."
    "Pick the odd word out: Spain, France, German, England, Singapore."
    "Spain is a country, France is a country, German is a language, ..."
    "So the odd one is German."
    "Pick the odd word out: {OPTIONS}"
    "[REASONING]"
    "[RESULT]"
from "EleutherAI/gpt-j-6B"
where
    not "Än" in REASONING and not "Pick" in REASONING and
    stops_at(REASONING, "Pick the odd word") and stops_at(REASONING, "Än") and
    stops_at(REASONING, "So the odd one") and stops_at(REASONING, ".") and len(WORDS(REASONING)) < 40
distribute
    RESULT over OPTIONS.split(", ")
```
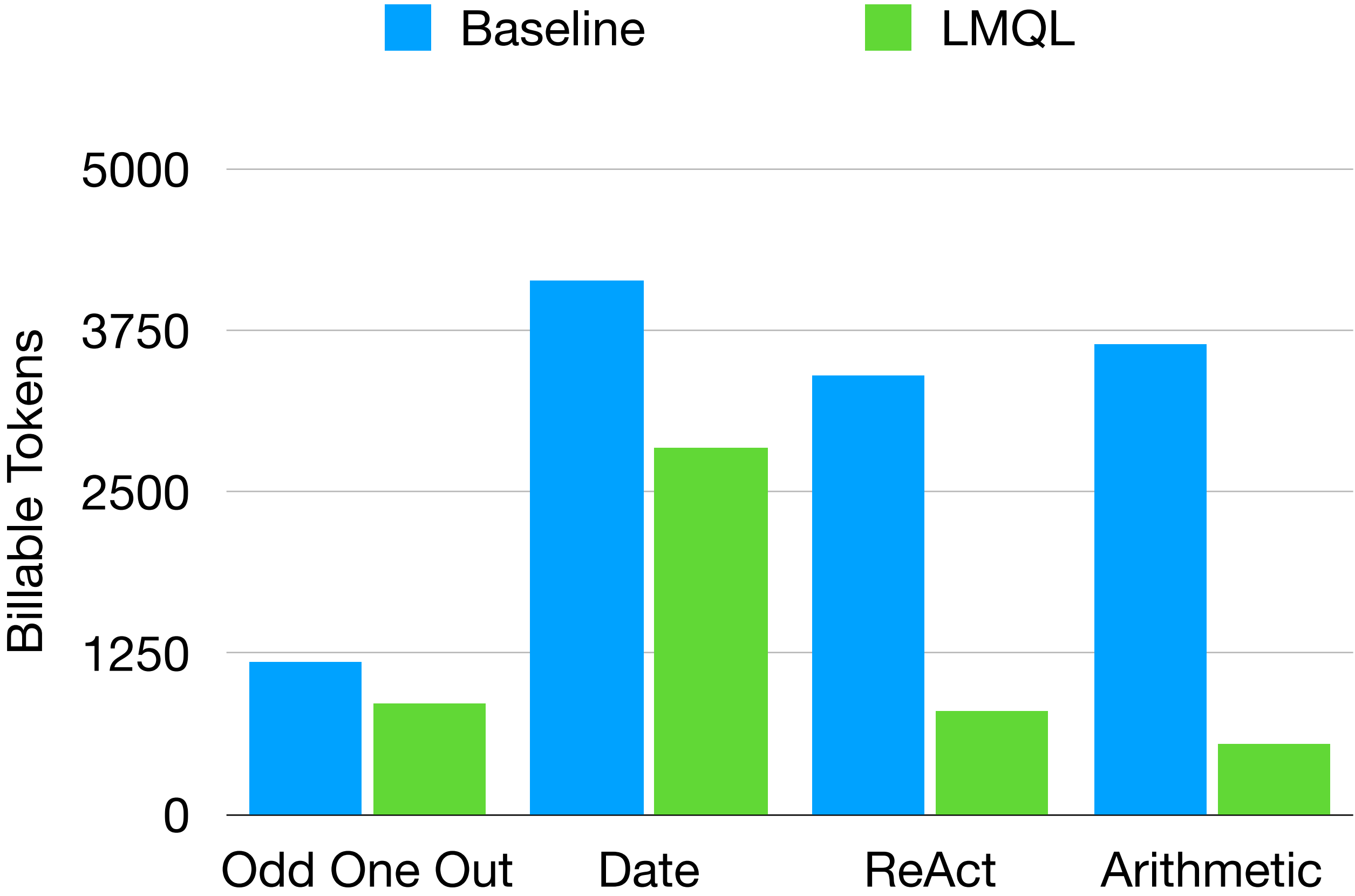
# Expressiveness

- Interactive Prompting

```
import wikipedia_utils
sample(no_repeat_ngram_size=3)
    "What is the elevation range for the area that the eastern sector of the Colorado orogeny extends into?"
    "Tho 1: I need to search Colorado orogeny, find the area that the eastern sector of the Colorado ..."
    "Act 2: Search 'Colorado orogeny'"
    "Obs 2: The Colorado orogeny was an episode of mountain building (an orogeny) ..."
    "Tho 3: It does not mention the eastern sector. So I need to look up eastern sector."
    ...
    "Tho 4: High Plains rise in elevation from around 1,800 to 7,000 ft, so the answer is 1,800 to 7,000 ft."
    "Act 5: Finish '1,800 to 7,000 ft'"
    "Where is Apple Computers headquartered?"
    for i in range(1024):
        "[MODE] {i}:"
    if MODE == "Tho":
        "[THOUGHT] "
    elif MODE == "Act":
        " [ACTION] '[SUBJECT]Än"
        if ACTION == "Search":
            result = wikipedia_utils.search(SUBJECT[:-1]) # cutting of the consumed '
            "Obs {i}: {result}Än"
        else:
            break # action must be FINISH
from "gpt2-xl"
where
    MODE in ["Tho", "Act"] and stops_at(THOUGHT, "Än") and
    ACTION in ["Search", "Finish"] and len(words(THOUGHT)) > 2 and
    stops_at(SUBJECT, "'") and not "Tho" in THOUGHT
```

# Expressiveness

- Arithmetic Reasoning

```
argmax(distribution_batch_size=1, max_length=2048)
    "(few-shot examples)"
    "Q: {QUESTION}Än"
    "A: Let's think step by step.Än"
    for i in range(1024):
        "[REASON_OR_CALC]"
        if REASON_OR_CALC.endswith("<<"):
            " [EXPR] "
            result = calculator.run(EXPR)
            " {result} >> "
        elif REASON_OR_CALC.endswith("So the answer"):
            break
        " is [RESULT]"
from "EleutherAI/gpt-j-6B"
where
    int(RESULT) and
    stops_at(REASON_OR_CALC, "<<") and
    stops_at(EXPR, "=") and
    stops_at(REASON_OR_CALC, "So the answer")
```

# Performance
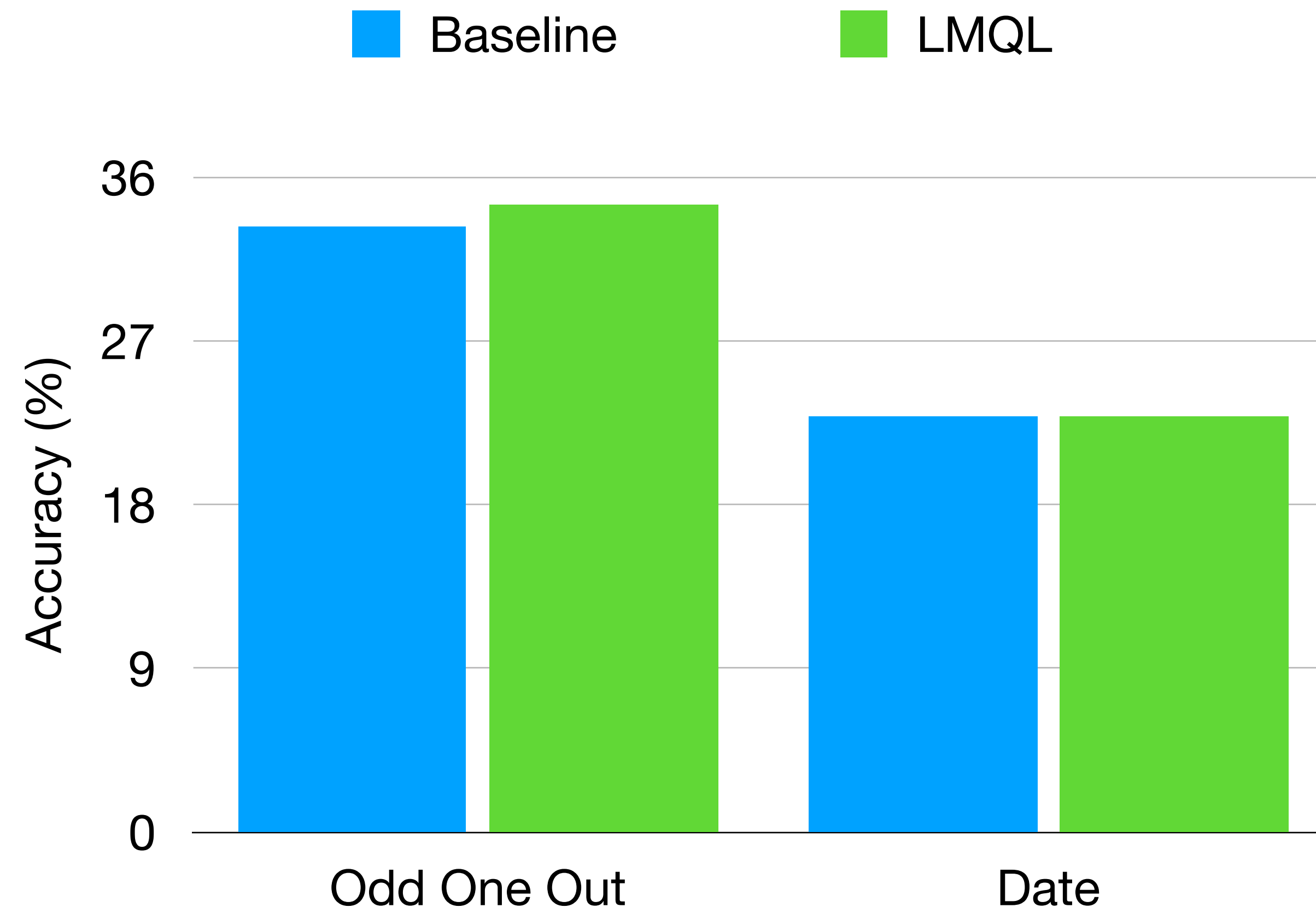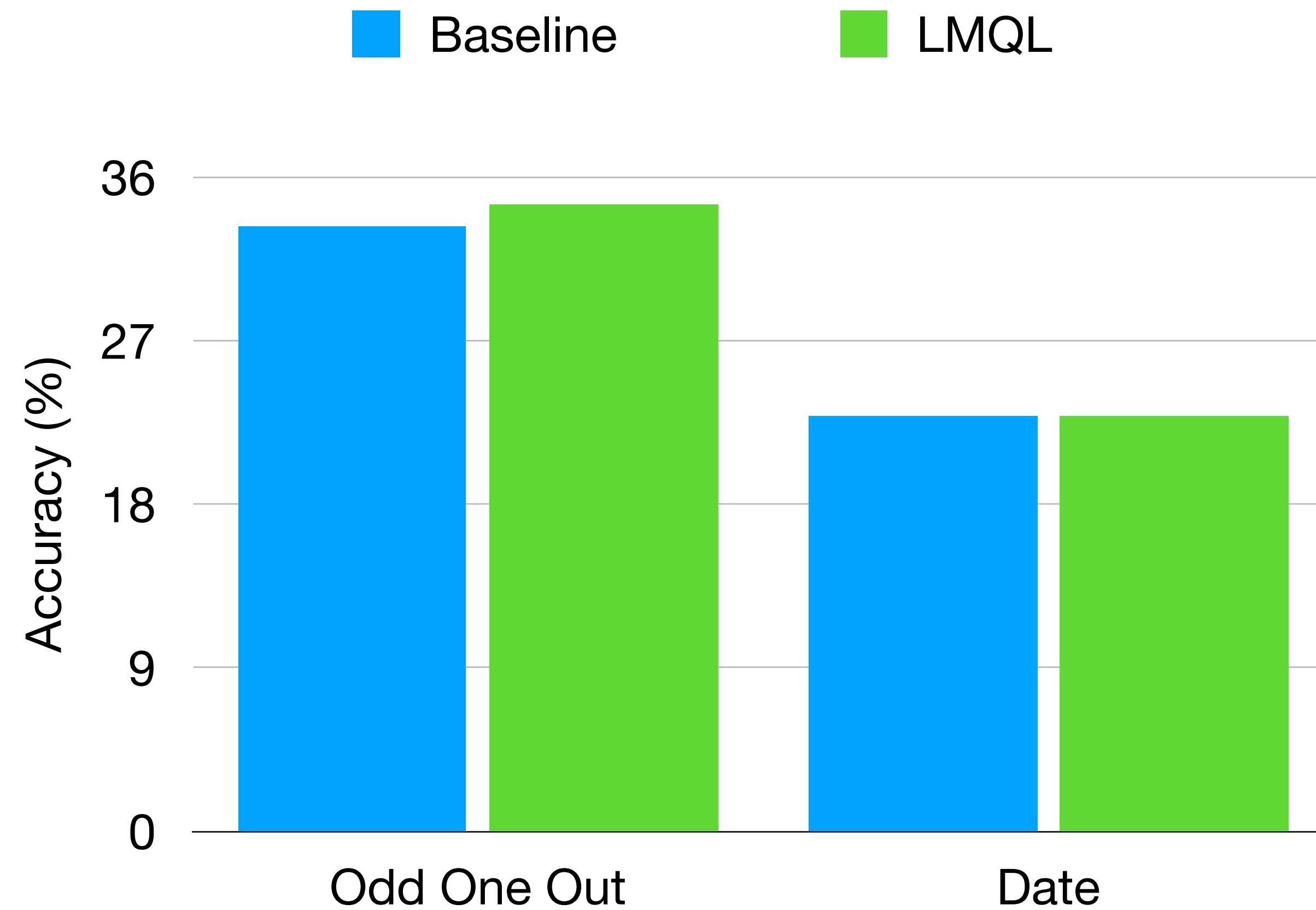
# Performance



Using much fewer tokens through control flow and variables

# Accuracy

# Accuracy



Accuracy isn't affected by LMQL

# Review

- Strong points

  - Formally defined constraint decoding method and soundness

  - Tools that can be useful in the wild

  - Serving nice visualization

- Weak points

  - Evaluated constraints are too easy

  - Not compositional (Nested Query)

## Query

```
beam(n=4)
    """English to French Translation:
    English: I am going to the store
    French: [TRANSLATION]
    """
from
    "openai/text-davinci-001"
where
    STOPS_AT(TRANSLATION, "\n")
```

▶ Run    ● Ready

Tokens: 1257, Req.: 16, Avb: 2.44
Time: 19.6s, 139.67 tok/step Est. Cost $0.0251

## Model Response
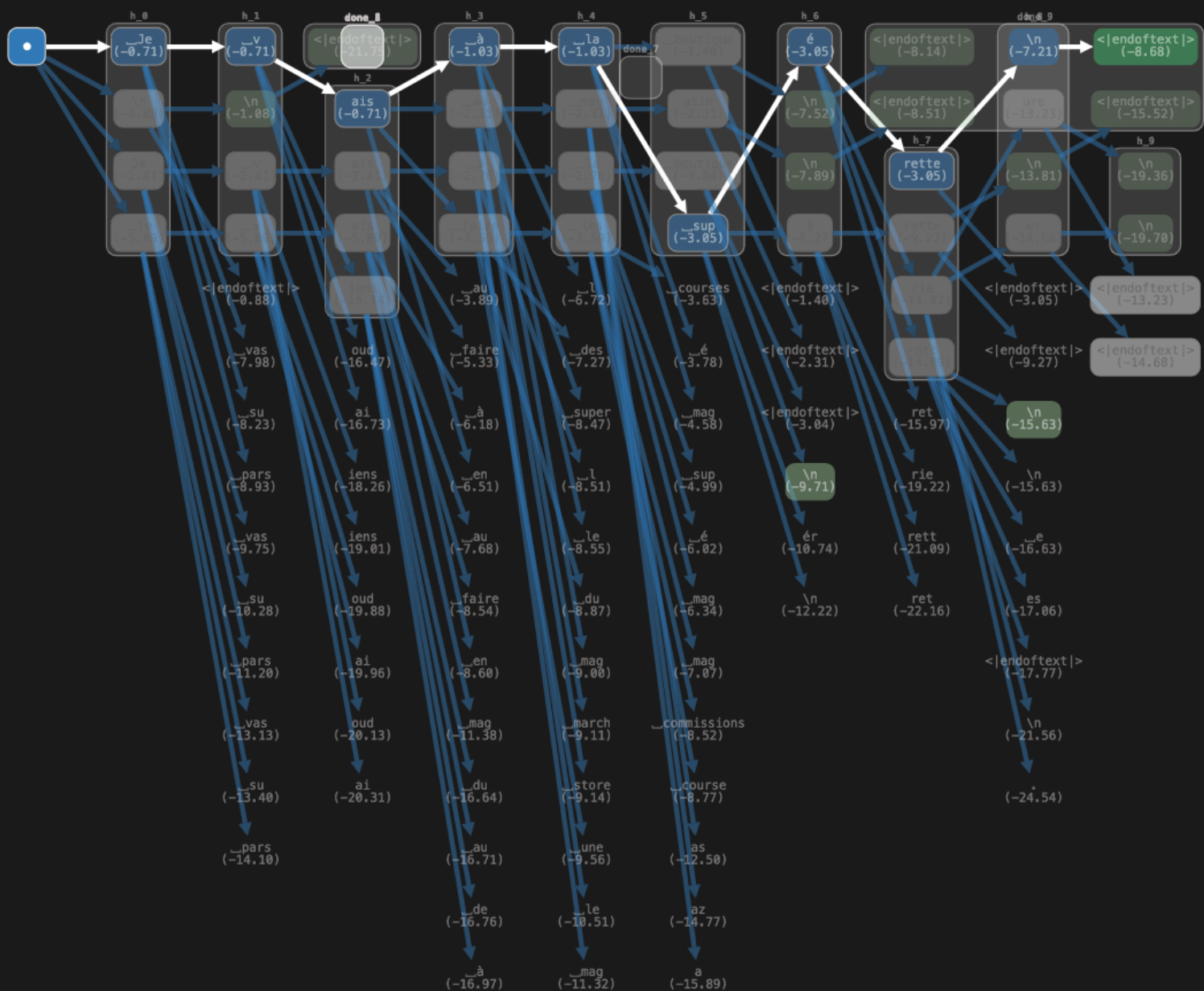
☑ Color By Variable    ☐ Show Latest

English to French Translation: ↵
English: I am going to the store ↵
French: [TRANSLATION] Je vais à la supérette <eos>

## Decoder Graph

☐ Fit    ☐ Eager Layouting



## Inspector

**DECODER**

LOGPROB          -1.4701520204544067
SEQLOGPROB       -8.682374383195718
POOL             "done_9"

**INTERPRETER**

VARIABLE         "__done__"
VALID            ● fin(true) ⊕
MASK             "<not available yet>"

**VARIABLES**

TRANSLATION      fin(" Je vais à la supérette")

**MISC**

DETERMINISTIC    true
DETERMINISTIC_5  [false,false,false,true,true]
DONE             true
NEXT_IDS         []
NEXT_LOGPROBS    []
PROMPT           "English to French Translation:\nEnglish: I am going to the
                 store\nFrench: Je vais à la supérette\n"
QUERY_HEAD       "<InterpretationHead query, (), {'STOPS_AT': None, 'TRANSLAT
                 ION': None, 'beam': None}>"

# Summary

- The paper suggests **a new paradigm** of interacting with the language model and its **implementation**

  - **LMP**: Language Model + Programming

  - **LMQL**: Add abstraction of language model to programming language

- Provide **control flow + variable + constraint decoding**

- We can use LMQL for **most of the existing prompting strategies**

- Lowering cost and slightly increasing accuracy

# Example2

```
argmax

  "A list of things not to forget
    when travelling:\n"

  things = []

  for _ in range(2):

    "- [THING]\n"

    things.append(THING)

from "gpt-3"

where

  len(words(THING)) <= 2
```

# Example2

```
argmax          ─────────▶   We are going to use greedy decoding

  "A list of things not to forget
   when travelling:\n"

  things = []

  for _ in range(2):

    "- [THING]\n"

    things.append(THING)

from "gpt-3"

where

  len(words(THING)) <= 2
```

# Example2

```
argmax

  "A list of things not to forget
   when travelling:\n"

  things = []

  for _ in range(2):

    "- [THING]\n"

    things.append(THING)

from "gpt-3"  ──────────▶  Use language model gpt2-large

where

  len(words(THING)) <= 2
```

45

# Example2

```
argmax

    "A list of things not to forget
      when travelling:\n"

    things = []

    for _ in range(2):

      "- [THING]\n"

      things.append(THING)

from "gpt-3"

where

    len(words(THING)) <= 2  ──────────→  THING is 2 words or less
```

# Example2

```
argmax

   "A list of things not to forget
    when travelling:\n"

   things = []  ──────────▶  Add variable things to environment

   for _ in range(2):

   "- [THING]\n"

   things.append(THING)

from "gpt-3"

where

   len(words(THING)) <= 2
```

"A list of …"

# Example2

```
argmax

  "A list of things not to forget
   when travelling:\n"

  things = []

  for _ in range(2):

    "- [THING]\n"  ──────────⟶  Call LM and Assign response to THING

    things.append(THING)

from "gpt-3"

where

  len(words(THING)) <= 2
```

**Environment**

things = []

**Prompt**

"A list of …
_ "

48

# Example2

```
argmax

    "A list of things not to forget
        when travelling:\n"

    things = []

    for _ in range(2):

        "- [THING]\n"

        things.append(THING)  ───────────▶  Append THING to things

from "gpt-3"

where

    len(words(THING)) <= 2
```

**Environment**

```
things = []
THING = "money"
```

**Prompt**

"A list of …
- money"

# Example2

```
argmax

   "A list of things not to forget
      when travelling:\n"

   things = []

   for _ in range(2):

      "- [THING]\n"                    Append THING to things

      things.append(THING)

from "gpt-3"

where

   len(words(THING)) <= 2
```

**Environment**

```
things =
["money"]
THING = "money"
```

**Prompt**

"A list of …
- money
- "

Append **THING** to **things**