

Language Model Programming Lecture 4: Prompting as Programming

Kyle Richardson, **Gijs Wijnholds**

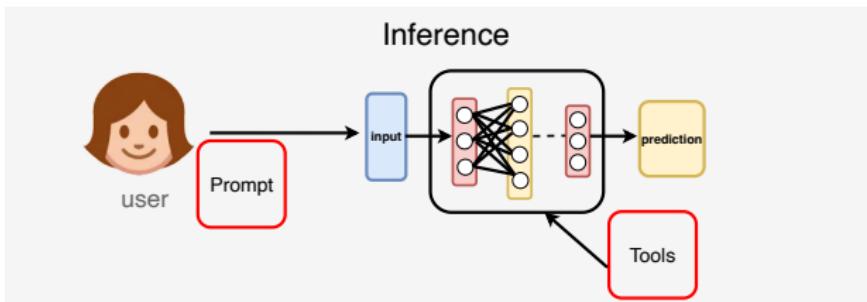
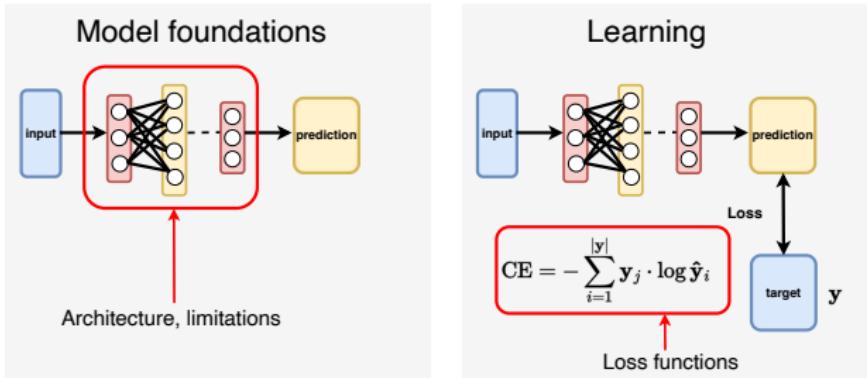
Allen Institute for AI (AI2)
Leiden Institute of Advanced Computer Science

August 2024

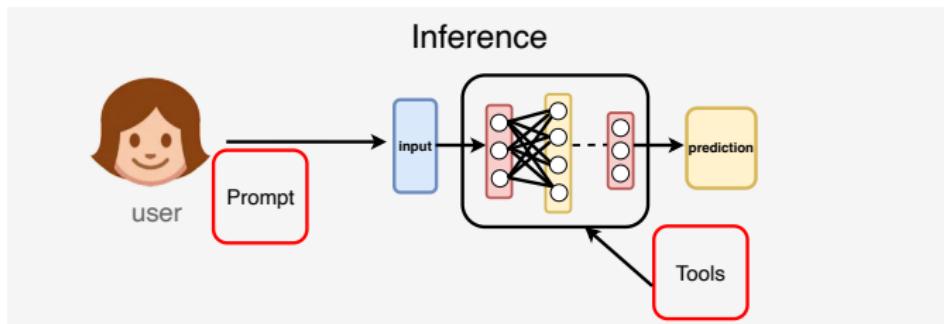


Universiteit
Leiden

Course Overview



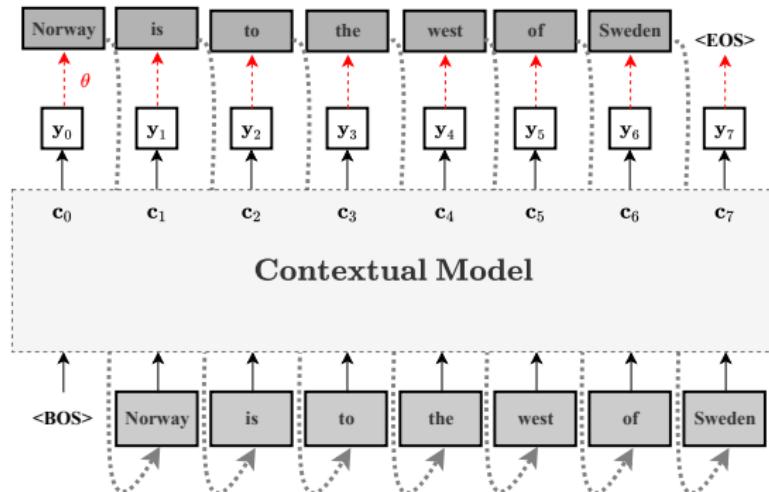
Today's Focus



- ▶ How to generate
- ▶ Prompting techniques
- ▶ Break
- ▶ Prompting is programming

Basics of generation with LLMs

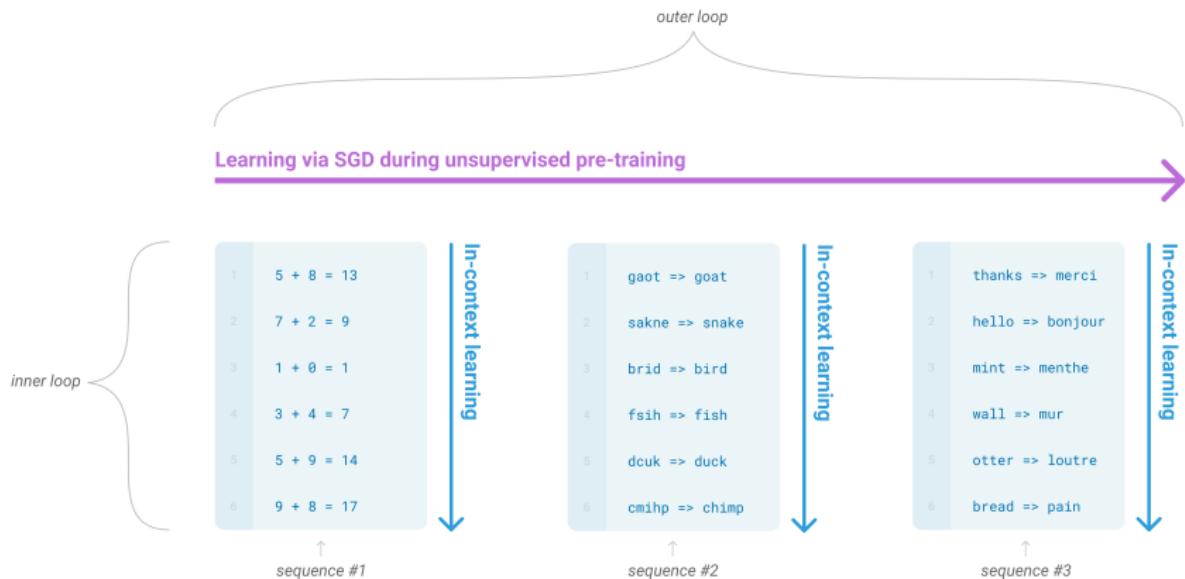
Prompting and decoding



- ▶ Learnt LLM probabilities are conditioned on the **prompt**
- ▶ The generated output is the combination of:
 - ▶ The prompt-conditioned probabilities
 - ▶ A **decoding strategy**

In-Context Learning

The power of scale during pretraining, the model learns strong “pattern recognition abilities”



Brown et al. [2020]

Prompting algorithm: simple decoding

Algorithm 1: Simple Generation

Input: LM f , tokenized prompt x

Output: Completion o

```
1:  $o \leftarrow \epsilon$ 
2: while  $True$  do
3:    $\mathbf{z} \leftarrow \text{softmax}(f(x \cdot o))$       // compute logits
4:    $t \leftarrow \text{pick}(\mathbf{z})$                 // decoding step
5:   if  $t = EOS$  then break
6:    $o \leftarrow o \cdot t$ 
7: end while
8: return  $o$ 
```

Prompting algorithm: simple decoding

Algorithm 2: Simple Generation

Input: LM f , tokenized prompt x

Output: Completion o

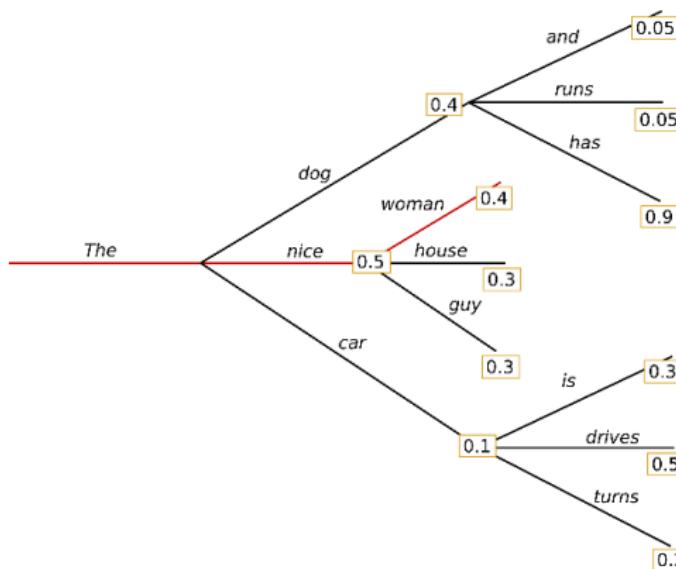
```
1:  $o \leftarrow \epsilon$ 
2: while  $True$  do
3:    $z \leftarrow \text{softmax}(f(x \cdot o))$       // compute logits
4:    $t \leftarrow \text{pick}(z)$                   // decoding step
5:   if  $t = EOS$  then break
6:    $o \leftarrow o \cdot t$ 
7: end while
8: return  $o$ 
```

So, how to pick?

Greedy Search

Strategy take the most probable option at each timestep:

$$w_t = \arg \max_w P(w \mid w_1 \dots w_{t-1})$$

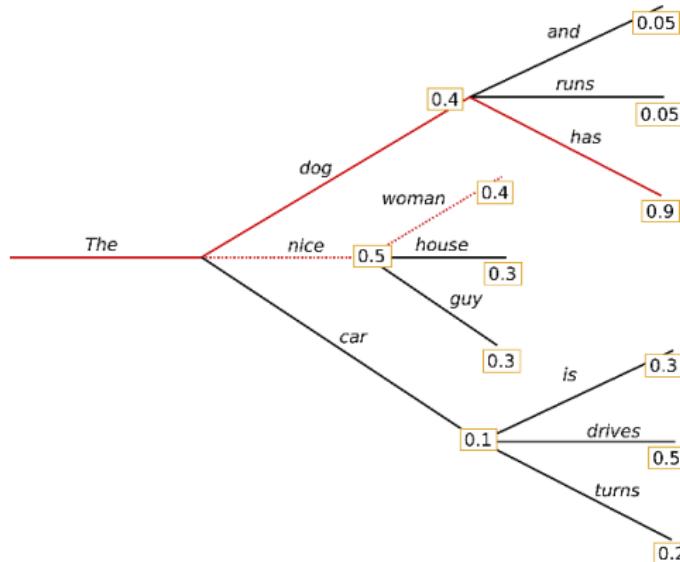


Pro linear in the sequence length

Con local choices don't guarantee the most probable sequence

Beam Search

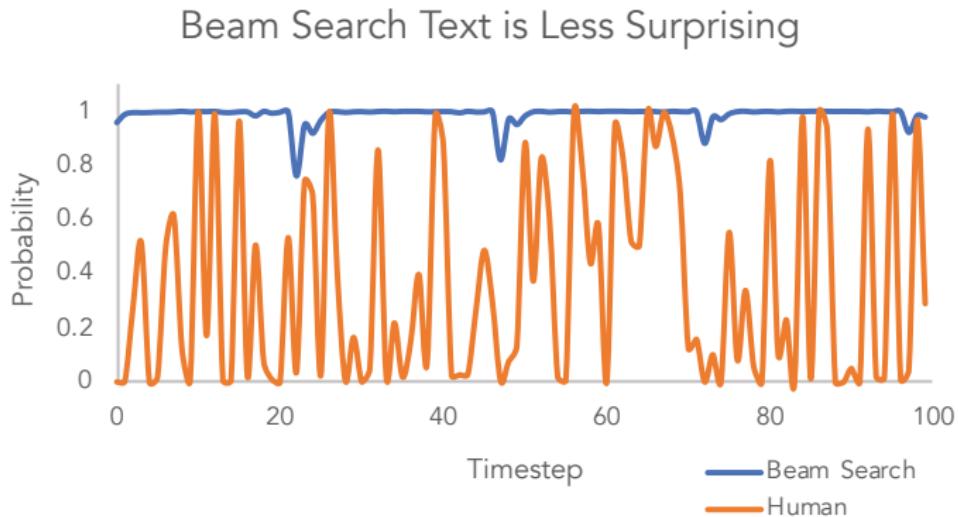
Strategy like greedy search, but keeping track of the k most likely beams



Pro more likely to find probable solutions

Why greedy/beam search isn't enough

Beam Search vs. Humans humans are much more varied than Beam Search:

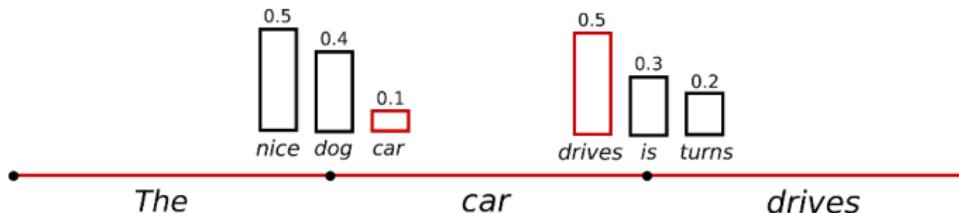


Holtzman et al. [2020]

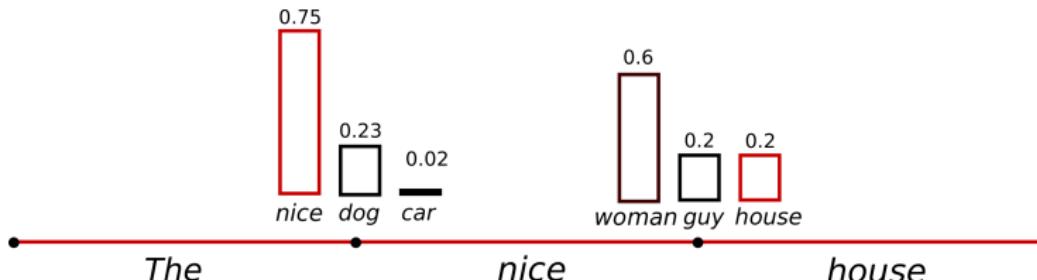
Sampling

Weighted sampling sample word from the conditional probability distribution:

$$w_t \sim P(w \mid w_1 \dots w_{t-1})$$

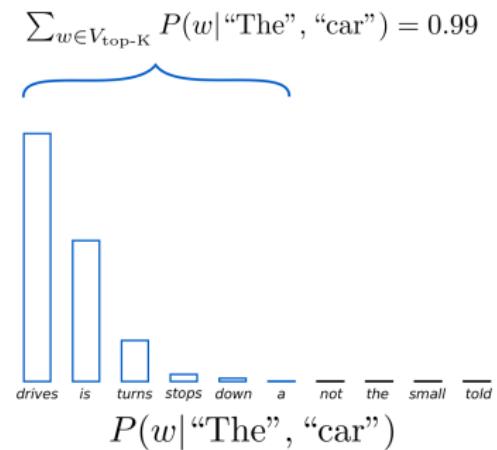
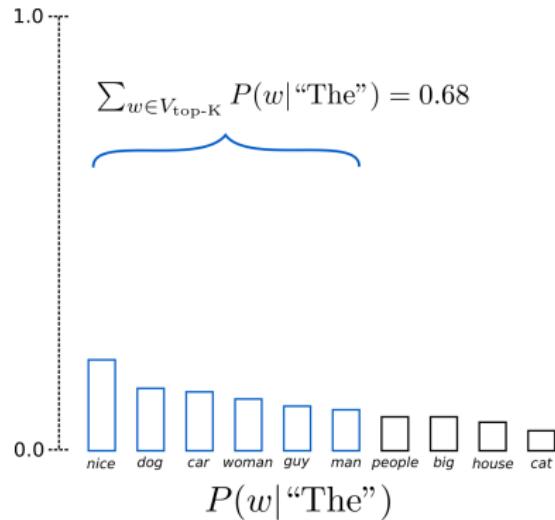


With temperature scaling



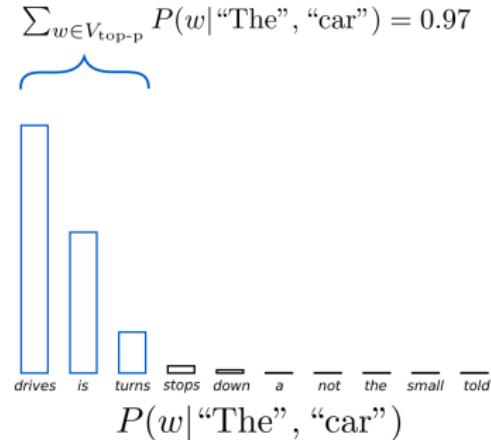
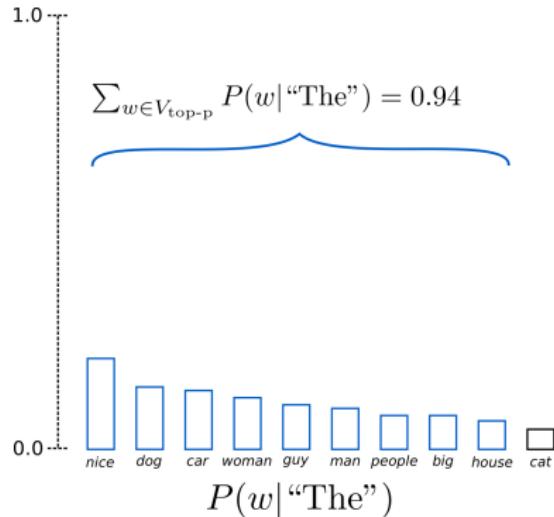
Top-K Sampling

Core idea sample only from the k most likely words:



Top-P Nucleus Sampling

Core idea sample from the set of most likely words whose combined probability exceeds p :



Prompting techniques

K-shot learning

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.



Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



Brown et al. [2020]

Chain of Thought Prompting

Standard Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27.

Chain-of-Thought Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had $23 - 20 = 3$. They bought 6 more apples, so they have $3 + 6 = 9$. The answer is 9.

Wei et al. [2022]

Chain of Thought Prompting: the Zero Shot case

(a) Few-shot

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) The answer is 8. X

(b) Few-shot-CoT

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) The juggler can juggle 16 balls. Half of the balls are golf balls. So there are $16 / 2 = 8$ golf balls. Half of the golf balls are blue. So there are $8 / 2 = 4$ blue golf balls. The answer is 4. ✓

(c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: The answer (arabic numerals) is

(Output) 8 X

(d) Zero-shot-CoT (Ours)

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: **Let's think step by step.**

(Output) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls. ✓

Least-to-most Prompting

Stage 1: Decompose Question into Subquestions

Q: It takes Amy 4 minutes to climb to the top of a slide. It takes her 1 minute to slide down. The water slide closes in 15 minutes. How many times can she slide before it closes?

Language Model

A: To solve "How many times can she slide before it closes?", we need to first solve: "How long does each trip take?"

Stage 2: Sequentially Solve Subquestions

It takes Amy 4 minutes to climb to the top of a slide. It takes her 1 minute to slide down. The slide closes in 15 minutes.

Subquestion 1 Q: How long does each trip take?

Language Model

A: It takes Amy 4 minutes to climb and 1 minute to slide down. $4 + 1 = 5$. So each trip takes 5 minutes.

Append model answer to Subquestion 1

It takes Amy 4 minutes to climb to the top of a slide. It takes her 1 minute to slide down. The slide closes in 15 minutes.

Q: How long does each trip take?

A: It takes Amy 4 minutes to climb and 1 minute to slide down. $4 + 1 = 5$. So each trip takes 5 minutes.

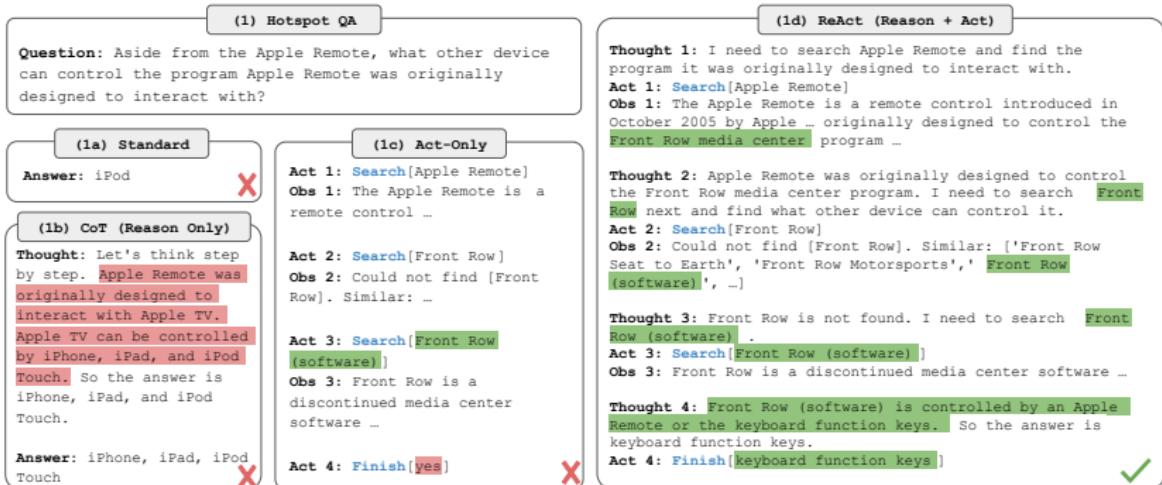
Subquestion 2

Q: How many times can she slide before it closes?

Language Model

A: The water slide closes in 15 minutes. Each trip takes 5 minutes. So Amy can slide $15 \div 5 = 3$ times before it closes.

ReAct prompting



Yao et al. [2023]

LM-generated program calls

Chain-of-Thought (Wei et al., 2022)

Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 tennis balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

Q: The bakers at the Beverly Hills Bakery baked 200 loaves of bread on Monday morning. They sold 93 loaves in the morning and 39 loaves in the afternoon. A grocery store returned 6 unsold loaves. How many loaves of bread did they have left?

Model Output

A: The bakers started with 200 loaves. They sold 93 in the morning and 39 in the afternoon. So they sold $93 + 39 = 132$ loaves. The grocery store returned 6 loaves. So they had $200 - 132 - 6 = 62$ loaves left.

The answer is 62.



Program-aided Language models (this work)

Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 tennis balls.

`tennis_balls = 5`

`2 cans of 3 tennis balls each is`

`bought_balls = 2 * 3`

`tennis_balls`. The answer is

`answer = tennis_balls + bought_balls`

Q: The bakers at the Beverly Hills Bakery baked 200 loaves of bread on Monday morning. They sold 93 loaves in the morning and 39 loaves in the afternoon. A grocery store returned 6 unsold loaves. How many loaves of bread did they have left?

Model Output

A: The bakers started with 200 loaves

`loaves_baked = 200`

`They sold 93 in the morning and 39 in the afternoon`

`loaves_sold_morning = 93`

`loaves_sold_afternoon = 39`

`The grocery store returned 6 loaves.`

`loaves_returned = 6`

The answer is

`answer = loaves_baked - loaves_sold_morning`
`- loaves_sold_afternoon + loaves_returned`

`>>> print(answer)`

74



Gao et al. [2023], Chen et al.

Interim Prompting Summary

- ▶ Including demonstration examples (few-shot reasoning)
- ▶ Inciting the model to generate reasoning steps
(Chain-of-Thought)
- ▶ Decomposing the question and solve with multiple prompts
(Least-to-most)
- ▶ Multiple steps including reasoning and actions (ReAct)
- ▶ Generating programs for reasoning (PAL)

Interim Prompting Summary

- ▶ Including demonstration examples (few-shot reasoning)
- ▶ Inciting the model to generate reasoning steps (Chain-of-Thought)
- ▶ Decomposing the question and solve with multiple prompts (Least-to-most)
- ▶ Multiple steps including reasoning and actions (ReAct)
- ▶ Generating programs for reasoning (PAL)

Let's embed (L)LMs inside
a programming language 😍

Prompting is programming

A query language for language models: LMQL

```
beam(n=3)
    "A list of good dad jokes. A indicates the "
    "punchline \n"
    "Q: How does a penguin build its house? \n"
    "A: Igloos it together. END \n"
    "Q: Which knight invented King Arthur's Round"
    "Table? \n"
    "A: Sir Cumference. END \n"
    "Q: [JOKE] \n"
    "A: [PUNCHLINE] \n"
from "gpt2-medium"
where
    STOPS_AT(JOKE, "?") and STOPS_AT(PUNCHLINE, "END")
    and len(words(JOKE)) < 20
    and len(characters(PUNCHLINE)) > 10
```

What is a programming language for prompting?

The requirements

- ▶ A description of the language (syntax)
- ▶ The interface to prompting (semantics)

The components

- ▶ Multiple variables in the prompt, that need to be filled
- ▶ Declarative constraints that generation needs to be constrained by

From query language to prompting

- ▶ Iterative generation of each variable,
- ▶ Token masking to implement constraints on decoding

The LMQL Language

LMQL Program

```
<decoder> <query>
from <model>
[where <cond>]
[distribute <dist>]
```

```
<decoder> ::= argmax | beam(n=<int>) | sample(n=<int>)
<query> ::= <python_statement>+
<cond> ::= <cond> and <cond> | <cond> or <cond> | not <cond> | <cond_term>
         | <cond_term> <cond_op> <cond_term>
<cond_term> ::= <python_expression>
<cond_op> ::= < | > | = | in
<dist> ::= <var> over <python_expression>
```

The query

Core idea a query is (informally) the body of a Python function

- ▶ without inner function declarations,
- ▶ where top-level strings are LM prompts

The query

Core idea a query is (informally) the body of a Python function

- ▶ without inner function declarations,
- ▶ where top-level strings are LM prompts

Prompt variables The prompts have two special components:

- ▶ "[varname]": a phrase that needs to be generated (a *hole*)
- ▶ "{varname)": a reference to a previously generated phrase

The query

Core idea a query is (informally) the body of a Python function

- ▶ without inner function declarations,
- ▶ where top-level strings are LM prompts

Prompt variables The prompts have two special components:

- ▶ "[varname]": a phrase that needs to be generated (a *hole*)
- ▶ "{varname)": a reference to a previously generated phrase

Evaluation of a query we iteratively evaluate top-level strings keeping track of:

- ▶ an *interaction trace* (all text generated so far)
- ▶ a *scope* (variable assignment)

Evaluation of an LMQL program

Algorithm 1: Evaluation of a top-level string s

Input: string s , trace u , scope σ , language model f

```
1 if  $s$  contains [ $\langle\text{varname}\rangle$ ] then
2    $s_{\text{pre}}, \text{varname}, s_{\text{post}} \leftarrow \text{unpack}(s)$ 
    // e.g. "a [b] c" → "a ", "b", " c"
3    $u \leftarrow us_{\text{pre}}$  // append to trace
4    $v \leftarrow \text{decode}(f, u)$  // use the LM for the hole
5    $\sigma[\text{varname}] \leftarrow v$  // updated scope
6    $u \leftarrow uv$  // append to trace
7 else if  $s$  contains { $\langle\text{varname}\rangle$ } then
8    $\text{varname} \leftarrow \text{unpack}(s)$  // e.g. "{b}" → "b"
9    $v \leftarrow \sigma[\text{varname}]$  // retrieve value from scope
10   $s \leftarrow \text{subs}(s, \text{varname}, v)$  // replace placeholder
    with value
11   $u \leftarrow us$  // append to trace
12 else
13  |  $u \leftarrow us$  // append to trace
14 end
```

No, really, prompting is programming

```
1 argmax
2   "A list of things not to forget when "
3   "travelling:\n"
4   things = []
5   for i in range(2):
6     "- [THING]\n"
7     things.append(THING)
8   "The most important of these is [ITEM]."
9 from "EleutherAI/gpt-j-6B"
10 where
11   THING in ["passport",
12           "phone",
13           "keys", ...] // a longer list
14   and len(words(ITEM)) <= 2
```

Evaluation example

update

$$u = \epsilon$$

$$g = \{\}$$

state after update

Evaluation example

update	state after update
$s \leftarrow "A_list_of_things_not_to_forget_when"$ $u \leftarrow us$	$u = "A_list_of_things_not_to_forget_when"$ $g = \{\}$

Evaluation example

update	state after update
$s \leftarrow "A_list_of_things_not_to_forget_when"$ $u \leftarrow us$	$u = \epsilon$ $g = \{\}$ $u = "A_list_of_things_not_to_forget_when"$ $g = \{\}$
$s \leftarrow "travelling:\n"$ $u \leftarrow us$	 $u = "A_list_of_things_not_to_forget_when\ travelling\n"$ $g = \{\}$

Evaluation example

update	state after update
$s \leftarrow "A_list_of_things_not_to_forget_when"$ $u \leftarrow us$	$u = \epsilon$ $g = \{\}$
$s \leftarrow "travelling:\n"$ $u \leftarrow us$	$u = "A_list_of_things_not_to_forget_when\ travelling,\n"$ $g = \{\}$
$s \leftarrow "-[THING]\n"$ $s_{pre}, varname, s_{post} \leftarrow "-", THING, \n$ $u \leftarrow us_{pre}$ $v \leftarrow "sun_screen" = decode(f, u)$ $u \leftarrow uvs_{post}$ $g[varname] \leftarrow v$	$u = "A_list_of_things_not_to_forget_when\ travelling,\n -sun_screen\n"$ $g = \{i = 0, THING = "sun_screen",$ $things = ["sun_screen"]\}$

Evaluation example

update	state after update
	$u = \epsilon$ $g = \{ \}$
$s \leftarrow "A_list_of_things_not_to_forget_when"$ $u \leftarrow us$	$u = "A_list_of_things_not_to_forget_when"$ $g = \{ \}$
$s \leftarrow "travelling:\n"$ $u \leftarrow us$	$u = "A_list_of_things_not_to_forget_when\ travelling\n"$ $g = \{ \}$
$s \leftarrow "-__["THING"]\n"$ $s_{pre}, varname, s_{post} \leftarrow "-__", THING, \n$ $u \leftarrow u s_{pre}$ $v \leftarrow "sun_screen" = decode(f, u)$ $u \leftarrow u v s_{post}$ $g[varname] \leftarrow v$	$u = "A_list_of_things_not_to_forget_when\ travelling\n -__sun_screen\n"$ $g = \{ i = 0, THING = "sun_screen",$ $\quad \quad \quad things = ["sun_screen"] \}$
$s \leftarrow "-__["THING"]\n"$ $s_{pre}, varname, s_{post} \leftarrow "-__", THING, \n$ $u \leftarrow u s_{pre}$ $v \leftarrow "beach_towel" = decode(f, u)$ $u \leftarrow u v s_{post}$ $g[varname] \leftarrow v$	$u = "A_list_of_things_not_to_forget_when\ travelling\n -__sun_screen\n -__beach_towel\n"$ $g = \{ i = 1, THING = "beach_towel",$ $\quad \quad \quad things = ["sun_screen", "beach_towel"] \}$

Distribution

Core idea grabbing the probability distribution of the LM for the last variable, in a given set

Distribution

Core idea grabbing the probability distribution of the LM for the last variable, in a given set

```
1 argmax
2 "A list of things not to forget when "
3 "travelling:\n"
4 things = []
5 for i in range(2):
6     "- [THING]\n"
7     things.append(THING)
8 "The most important of these is [ITEM]."
9 from "EleutherAI/gpt-j-6B"
10 where
11     THING in ["passport",
12                 "phone",
13                 "keys", ...] // a longer list
14 distribute ITEM over things
```



A list of things not to forget when travelling:
- sun screen
- beach towel

The most important of these is {sun screen 65%
beach towel 35%

Constraints

Constraint syntax

```
<cond> ::= <cond> and <cond> | <cond> or <cond> | not <cond> | <cond_term>
          | <cond_term> <cond_op> <cond_term>
<cond_term> ::= <python_expression>
<cond_op> ::= < | > | = | in
```

*deterministic pure python functions + constants

Built-in functions

```
[w1, ..., wk] ← words(<var>)      //splits <var> into words w1, ..., wk
[s1, ..., sk] ← sentences(<var>)    //splits <var> into sentences s1, ..., sk
b ← stop_at(<var>, t)                  //indicates if <var> ends in token or string t
```

Constraints

Constraint syntax

```
<cond> ::= <cond> and <cond> | <cond> or <cond> | not <cond> | <cond_term>
          | <cond_term> <cond_op> <cond_term>
<cond_term> ::= <python_expression>
<cond_op> ::= < | > | = | in
```

*deterministic pure python functions + constants

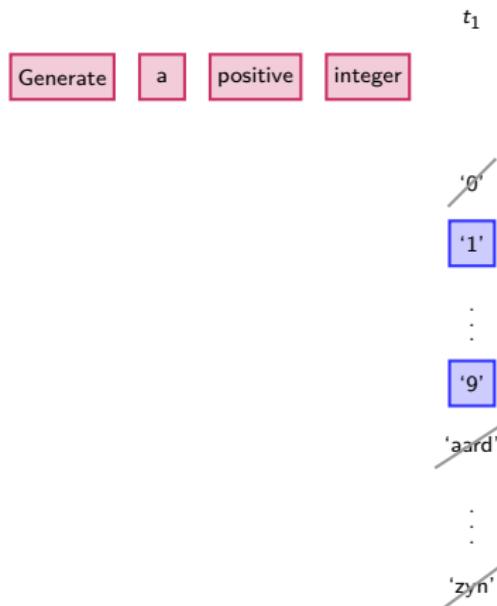
Built-in functions

```
[w1, ..., wk] ← words(<var>)      //splits <var> into words w1, ..., wk
[s1, ..., sk] ← sentences(<var>)    //splits <var> into sentences s1, ..., sk
b ← stop_at(<var>, t)                  //indicates if <var> ends in token or string t
```

Constraint validation how to evaluate expressions over to-be-generated strings?

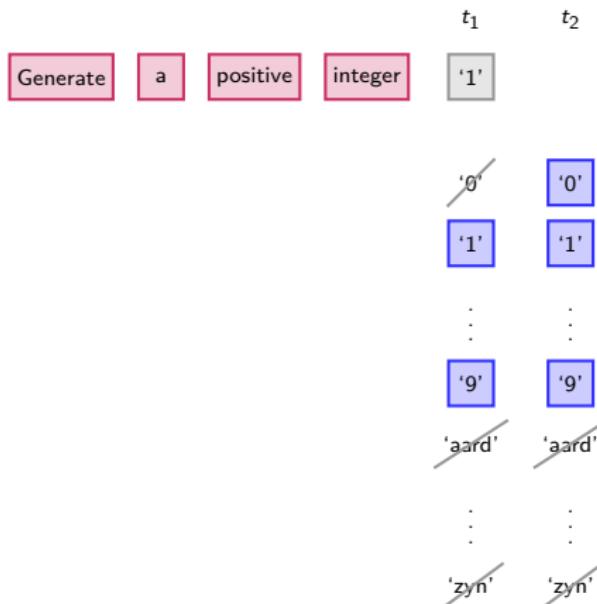
Masked Decoding

Core idea block out part of the vocabulary during decoding



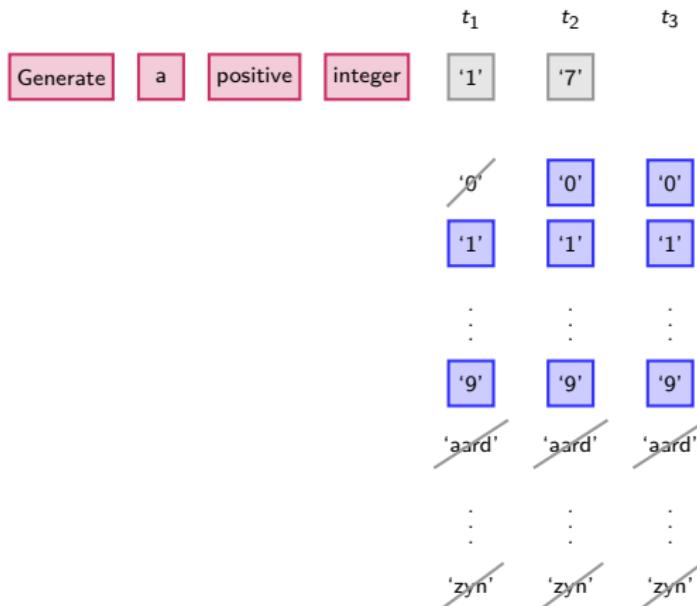
Masked Decoding

Core idea block out part of the vocabulary during decoding



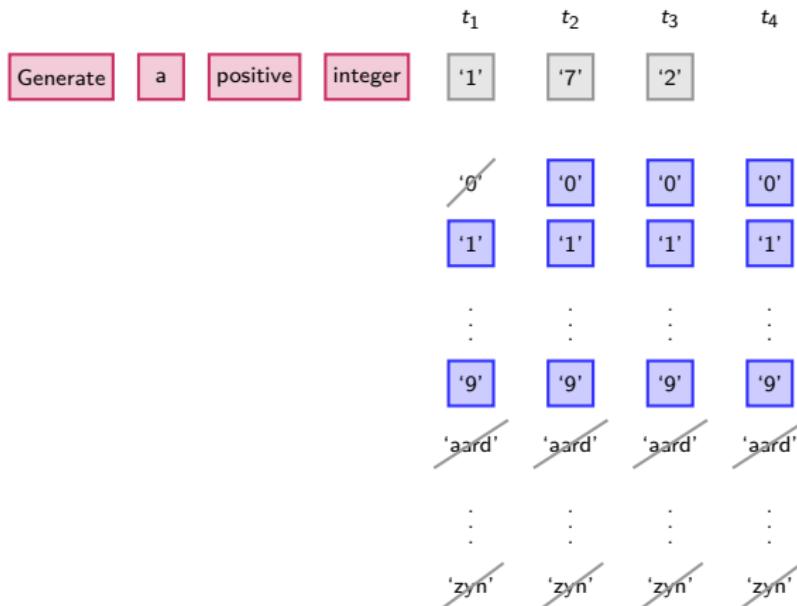
Masked Decoding

Core idea block out part of the vocabulary during decoding



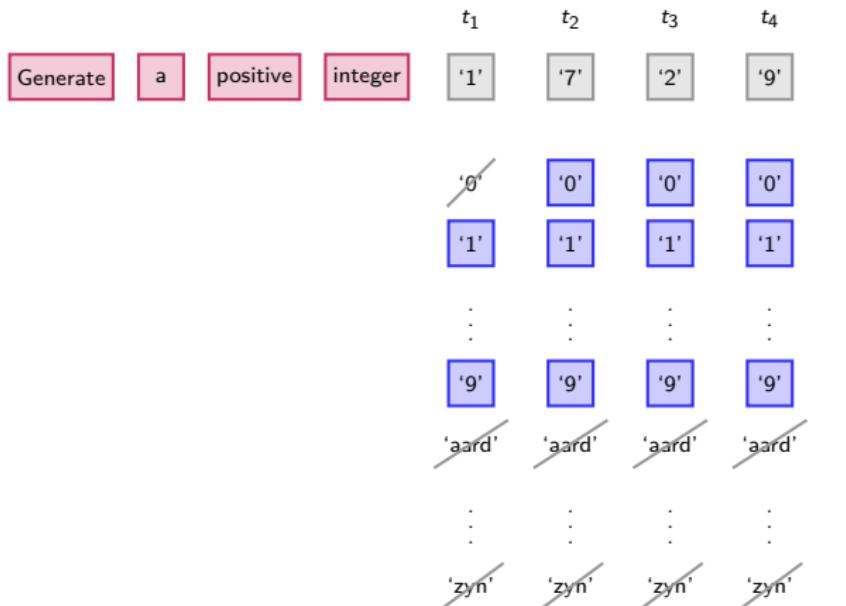
Masked Decoding

Core idea block out part of the vocabulary during decoding



Masked Decoding

Core idea block out part of the vocabulary during decoding



Constrained decoding with masks

Algorithm 2: Decoding

Input: trace u , scope σ , LM f

Output: decoded sequence v

```
1  $v \leftarrow \epsilon$ 
2 while True do
3    $m \leftarrow \text{compute\_mask}(u, \sigma, v)$ 
4   if  $\wedge_i(m_i = 0)$  then break
5    $z \leftarrow {}^1/z \cdot m \odot \text{softmax}(f(uv))$ 
6    $t \leftarrow \text{pick}(z)$ 
7   if  $t = \text{EOS}$  then break
8    $v \leftarrow vt$ 
9 end
```

Constrained decoding with masks

Algorithm 2: Decoding

Input: trace u , scope σ , LM f

Output: decoded sequence v

```
1  $v \leftarrow \epsilon$ 
2 while True do
3    $m \leftarrow \text{compute\_mask}(u, \sigma, v)$ 
4   if  $\wedge_i(m_i = 0)$  then break
5    $z \leftarrow {}^1/z \cdot m \odot \text{softmax}(f(uv))$ 
6    $t \leftarrow \text{pick}(z)$ 
7   if  $t = \text{EOS}$  then break
8    $v \leftarrow vt$ 
9 end
```

Why not just check after generation?

Beurer-Kellner et al. [2023]

Naive constrained decoding

Algorithm 3: Naive Decoding with Constraints

Input: trace u , scope σ , language model f
Output: decoded sequence v

1 **Function** $decode_step(f, u, v)$
2 $z \leftarrow \text{softmax}(f(uv))$
3 $\mathbf{m} \leftarrow \mathbf{1}^{|\mathcal{V}|}$
4 **do**
5 $t \leftarrow \text{pick}(^1/z \cdot \mathbf{m} \odot z)$
6 **if** $t \neq EOS$ **then** $decode_step(f, u, vt)$
7 **else if** $t = EOS \wedge check(u, vt)$ **then**
8 **return** v
9 **else** $\mathbf{m}[t] \leftarrow 0$
9 **while** $\bigvee_i m_i = 1$
10 **decode_step**(f, u, ϵ)

Incrementalizing constraints

Requirement we need to compile constraints into a set of invalid tokens (the token mask)

Incrementalizing constraints

Requirement we need to compile constraints into a set of invalid tokens (the token mask)

Partial evaluation to model constraint satisfaction on sequences that are not yet fully generated

Incrementalizing constraints

Requirement we need to compile constraints into a set of invalid tokens (the token mask)

Partial evaluation to model constraint satisfaction on sequences that are not yet fully generated

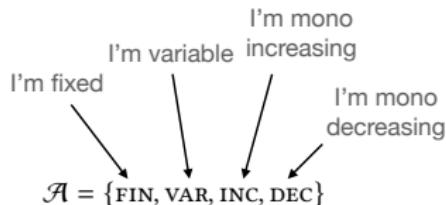
Final and Follow

- ▶ Final semantics: does an *expression* evaluate to a fixed value or not?
- ▶ Follow semantics: given an *expression* e , an *interaction trace* u and a *token* t , is ut valid?

Final Semantics

Values & Annotators

This is
my value
 $\llbracket e \rrbracket_\sigma = v$



Base case

expression **FINAL**[\cdot ; σ]

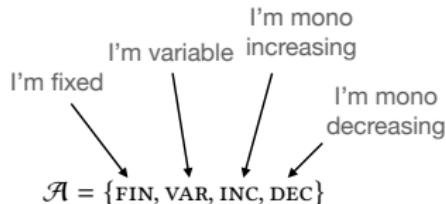
$\langle \text{const} \rangle$	FIN
python variable $\langle \text{pyvar} \rangle$	VAR
previous hole $\langle \text{var} \rangle$	FIN
current var $\langle \text{var} \rangle$	INC
future hole $\langle \text{var} \rangle$	INC

$\text{words}(v)$	$\text{FINAL}[v]$
$\text{sentences}(v)$	$\text{FINAL}[v]$
$\text{len}(v)$	$\text{FINAL}[v]$

Final Semantics

Values & Annotators

This is my value
 $\llbracket e \rrbracket_\sigma = v$



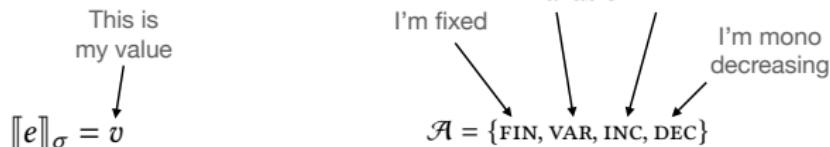
Equality & Functions

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

number equality $n == m$	$\begin{cases} \text{FIN} & \text{if } \text{FINAL}[n] = \text{FIN} \\ & \wedge \text{FINAL}[m] = \text{FIN} \\ \text{VAR} & \text{else} \end{cases}$	$\text{FINAL}[\text{THING} == \text{ITEM}]$	FIN	INC
string equality $x == y$	$\begin{cases} \text{FIN} & \text{if } \text{FINAL}[x] = \text{FIN} \\ & \wedge \text{FINAL}[y] = \text{FIN} \\ \text{VAR} & \exists i \bullet x[i] \neq y[i] \\ & \wedge \text{FINAL}[x] \neq \text{VAR} \\ & \wedge \text{FINAL}[y] \neq \text{VAR} \end{cases}$	$"\text{passport"} == \text{"pass"}$	VAR	FIN
function $\text{fn}(\tau_1, \dots, \tau_k)$	$\begin{cases} \text{FIN} & \text{if } \bigwedge_{i=1}^k \text{FINAL}[\tau_i] = \text{FIN} \\ \text{VAR} & \text{else} \end{cases}$	$"\text{passport"} == \text{"password"}$		

Final Semantics

Values & Annotators



Stop_at & Membership

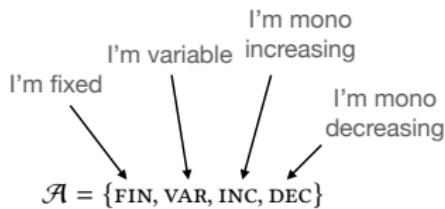
expression	$\text{FINAL}[\cdot ; \sigma]$
<code>stop_at(var, s)</code>	$\begin{cases} \text{FIN} & \text{if } \llbracket \text{var} \rrbracket_\sigma \cdot \text{endswith}(s) \\ & \wedge \text{FINAL}[\text{var}] = \text{INC} \\ \text{VAR} & \text{else} \end{cases}$
$x \text{ in } s$ for strings x, s	$\begin{cases} \text{FIN} & \text{if } x \text{ in } s \wedge \text{FINAL}[x] = \text{FIN} \\ & \wedge \text{FINAL}[s] = \text{INC} \\ \text{VAR} & \text{else} \end{cases}$
$e \text{ in } l$ for string e , set l	$\begin{cases} \text{FIN} & \text{if } \nexists i \in l \bullet i \cdot \text{startswith}(e) \\ & \wedge \text{FINAL}[x] \in \{\text{INC, FIN}\} \\ & \wedge \text{FINAL}[l] = \text{FIN} \\ \text{VAR} & \text{else} \end{cases}$

Final Semantics

Values & Annotators

This is
my value

$$[e]_{\sigma} = v$$



Inequality & Booleans

expression FINAL[\cdot ; σ]

$x < y$	$\begin{cases} \text{FIN} & \text{if } x < y \wedge \text{FINAL}[x] \in \{\text{DEC}, \text{FIN}\} \\ & \wedge \text{FINAL}[y] \in \{\text{INC}, \text{FIN}\} \\ \text{VAR} & \text{else} \end{cases}$
$a \text{ and } b$	$\begin{cases} \text{FIN} & \text{if } \exists v \in \{a, b\} \bullet [v]_{\sigma}^F = \text{FIN}(\perp) \\ \text{FIN} & \text{if } \forall v \in \{a, b\} \bullet [v]_{\sigma}^F = \text{FIN}(\top) \\ \text{VAR} & \text{else} \end{cases}$
$a \text{ or } b$	$\begin{cases} \text{FIN} & \text{if } \exists v \in \{a, b\} \bullet [v]_{\sigma}^F = \text{FIN}(\top) \\ \text{FIN} & \text{if } \forall v \in \{a, b\} \bullet [v]_{\sigma}^F = \text{FIN}(\perp) \\ \text{VAR} & \text{else} \end{cases}$
$\text{not } a$	$\text{FINAL}[a]$

Follow Maps

Core idea Can we exclude token t for a given expression, and interaction trace u ?

Follow Maps

Core idea Can we exclude token t for a given expression, and interaction trace u ?

Some example definitions

expression	$\text{FOLLOW}[\cdot](u, t)$
$\langle \text{const} \rangle$	$\llbracket \langle \text{const} \rangle \rrbracket_\sigma$
python variable	$\llbracket \text{pyvar} \rrbracket_{\sigma[v \leftarrow vt]}$
$\langle \text{pyvar} \rangle$	$\llbracket \langle \text{var} \rangle \rrbracket_\sigma$
previous hole $\langle \text{var} \rangle$	$\begin{cases} \text{FIN}(v) & \text{if } t = \text{EOS} \\ \text{INC}(vt) & \text{else} \end{cases}$
current var v	$\begin{cases} \text{FIN}(v) & \text{if } t = \text{EOS} \\ \text{INC}(vt) & \text{else} \end{cases}$
future hole $\langle \text{var} \rangle$	None
$\text{sentences}(v)$	$\begin{cases} \text{FIN}(s_1, \dots, s_k) & \text{if } t = \text{EOS} \\ \text{INC}(s_1, \dots, s_k, t) & \text{if } s_k \cdot \text{endswith}(\text{"."}) \\ \text{INC}(s_1, \dots, s_k t) & \text{else} \end{cases}$ where $s_1, \dots, s_k \leftarrow \llbracket \text{sentences}(v) \rrbracket_\sigma$

Follow Maps

Core idea Can we exclude token t for a given expression, and interaction trace u ?

Membership

$$\begin{array}{ll} \textbf{expression} & \textbf{FOLLOW}[\cdot](u, t) \\ \\ \begin{matrix} x \text{ in } l \\ \text{for constant list/set } l \end{matrix} & \begin{cases} \text{FIN}(\top) & \text{if } t \text{ in } l \\ \text{VAR}(\perp) & \text{if } \exists e \in l \bullet \\ & \quad e.\text{startswith}(vt) \\ \text{FIN}(\perp) & \text{else} \end{cases} \end{array}$$

Hawking

$$\text{FOLLOW}[\text{TEXT in ["Stephen Hawking"]}]("Steph", "en") = \text{VAR}(\perp)$$

$$\text{FOLLOW}[\text{TEXT in ["Stephen Hawking"]}]("Steph", "anie") = \text{FIN}(\perp)$$

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

Sound approximation

Token maps using the Follow semantics we can decide on tokens that will definitely invalidate the constraints

Soundness

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

In other words we are guaranteed to never mask out valid tokens, but will catch as many invalid tokens as possible

Sound approximation

Token maps using the Follow semantics we can decide on tokens that will definitely invalidate the constraints

Soundness

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

In other words we are guaranteed to never mask out valid tokens, but will catch as many invalid tokens as possible

But how well does it actually work?

Experimental Evaluation

Chain of Thought as an LQML program

```
argmax
    "Pick the odd word out: skirt, dress, pen, jacket.\n"
    "skirt is clothing, dress is clothing, pen is an object, jacket is clothing.\n"
    "So the odd one is pen.\n\n"
    "Pick the odd word out: Spain, France, German, England, Singapore.\n"
    "Spain is a country, France is a country, German is a language, ... \n"
    "So the odd one is German.\n\n"
    "Pick the odd word out: {OPTIONS}\n"
    "[REASONING]"
    "[RESULT]"

from "EleutherAI/gpt-j-6B"
where
    not "\n" in REASONING and not "Pick" in REASONING and
    stops_at(REASONING, "Pick the odd word") and stops_at(REASONING, "\n") and
    stops_at(REASONING, "So the odd one") and stops_at(REASONING, ".") and
    len(WORDS(REASONING)) < 40
distribute
    RESULT over OPTIONS.split(", ")
```

ReAct as an LQML program

```
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 ...\\n"
"Act 2: Search 'Colorado orogeny'\\n"
"Obs 2: The Colorado orogeny was an episode of mountain building (an orogeny) ...\\n"
Tho 3: It does not mention the eastern sector. So I need to look up eastern sector.\\n"

...
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?\\n"
for i in range(1024):
    "[MODE] {i}:""
    if MODE == "Tho":
        "[THOUGHT] "
    elif MODE == "Act":
        "[ACTION] '[SUBJECT]\\n"
        if ACTION == "Search":
            result = wikipedia_utils.search(SUBJECT[:-1]) # cutting of the consumed '
            "Obs {i}: {result}\\n"
        else:
            break # action must be FINISH
from "gpt2-xl"
where
    MODE in ["Tho", "Act"] and stops_at(THOUGHT, "\\n") and
    ACTION in ["Search", "Finish"] and len(words(THOUGHT)) > 2 and
    stops_at(SUBJECT, "") and not "Tho" in THOUGHT
```

Experimental Comparison

	GPT-J-6B [27]				OPT-30B [34]			
	Standard Decoding	LMQL	Δ	Est. Cost Savings	Standard Decoding	LMQL	Δ	Est. Cost Savings
<i>Odd One Out</i>								
Accuracy	33.33%	34.52%	1.19%		34.52%	34.52%	0.00%	
Decoder Calls	7.96	5.96	-25.11%		7.96	5.96	-25.11%	
Model Queries	73.04	41.51	-43.16%		73.04	40.70	-44.27%	
Billable Tokens	1178.71	861.32	-26.93%	0.63¢/query	1173.21	856.17	-27.02%	0.63¢/query

Beurer-Kellner et al. [2023]

Conclusion

Prompting is programming we can embed LM calls inside of a programming language

Constraint validation we need to translate constraints into a token masking strategy

Sound approximation The translation in LMQL is sound but not complete, i.e. we don't mask out all invalid tokens

The cliffhanger

But what if we want to generate

The cliffhanger

But what if we want to generate
structured output 😍?

The cliffhanger

But what if we want to generate
structured output 😍?

► JSON 😍?

The cliffhanger

But what if we want to generate
structured output 😍?

- ▶ JSON 😍?
- ▶ Python 😍 😍?

The cliffhanger

But what if we want to generate
structured output 😍?

- ▶ JSON 😍?
- ▶ Python 😍 😍?
- ▶ First-order logic 😍 😍 😍?

The cliffhanger

But what if we want to generate
structured output 😍?

- ▶ JSON 😍?
- ▶ Python 😍 😍?
- ▶ First-order logic 😍 😍 😍?

Tomorrow!

Thank you.

References |

- Luca Beurer-Kellner, Marc Fischer, and Martin Vechev. Prompting is programming: A query language for large language models. *Proceedings of the ACM on Programming Languages*, 7(PLDI):1946–1969, 2023.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on Machine Learning Research*.
- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models. In *International Conference on Machine Learning*, pages 10764–10799. PMLR, 2023.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rygGQyrFvH>.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022.

References II

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35: 24824–24837, 2022.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=WE_vluYUL-X.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, et al. Least-to-most prompting enables complex reasoning in large language models. In *The Eleventh International Conference on Learning Representations*, 2023.