Ask, and it shall be given: On the Turing completeness of prompting Ruizhong Qiu, Zhe Xu, Wenxuan Bao, Hanghang Tong {rq5,zhexu3,wbao4,htong}@illinois.edu

HIGHLIGHTS

> Our work: the first theory on the LLM prompting paradigm.

- Existing theories: the *one-model-one-task* paradigm;
- C.f. LLM prompting: the one-model-many-tasks paradigm.

> Expressive power: Prompting is Turing-complete.

- Not only existence: a simple & explicit construction.
- \succ **Complexity bounds**: For any TIME(t(n)) function,
- **CoT complexity**: $O(t(n) \log t(n))$ CoT steps;
- **Precision complexity**: O(log(n + t(n))) bits.
- Nearly as efficient as the class of all Transformers.
- (Not covered in this poster. See our paper for detail...)

PRELIMINARIES

> Background:

- **Decoder-only Transformers:** mainstream architecture of LLMs.
- *Expressive power*: what functions a model class can represent.
- > Existing theories: classic one-model-one-task paradigm.
- The class of all hard-attention Transformers is Turing-complete;
- Any TIME(t(n)) function is computable in O(t(n)) CoT steps.

> Practice: LLM prompting (i.e., one-model-many-tasks).

- A single general-purpose LLM; different prompts for different tasks.
- Fundamentally, how powerful is this paradigm?

> Notations:

- Σ : an **alphabet** (i.e., the set of tokens).
- Σ^* : the set of strings (possibly empty) over Σ .
- Σ^+ : the set of non-empty strings over Σ .
- Σ^{ω} : the set of countably infinite strings over Σ .
- $\Gamma: \Sigma^+ \to \Sigma$: a **decoder-only** Transformer.
 - Γ predicts the next-token via **greedy** decoding.
- generate $_{\Gamma}: \Sigma^+ \to \Sigma^+ \cup \Sigma^{\omega}$: autoregressive generation using Γ .
- Here, the output of generate_{Γ} does not include the prompt.



- > Theorem: There exist
- - generate $\Gamma(\pi_{\varphi} \cdot \text{tokenize}(x))$ computes a finite CoT, and

readout (generate Γ (π)

Remarks:

- Σ, Γ , tokenize, and readout are independent of the function φ ;
- The prompt π_{φ} is independent of the input x;
- The rest of this poster is the proof sketch of this theorem...

TWO-TAPE POST-TURING MACHINES

> Key idea: Define a new imperative model of computation. **Setting**: two bi-infinite tapes A & B. Each tape has Infinitely many cells over the binary alphabet {0,1} and • A head pointing to a cell. Input is initially written to tape A. **Proposed model**: two-tape Post–Turing machines (2-PTMs). • A 2-PTM is defined by a finite instruction sequence $\iota = \langle \iota_0, ..., \iota_{|\iota|-1} \rangle$. Each instruction ι_i ($0 \le j < |\iota|$) is one of the following: (below, $\tau \in \{A, B\}$)

- #: halt;

 τ L: move the head of tape τ one cell left, and go to ι_{j+1} ; τ R: move the head of tape τ one cell right, and go to ι_{j+1} ; • $\tau 0$: write 0 to the pointed cell of tape τ , and go to ι_{j+1} ; τ 1: write 1 to the pointed cell of tape τ , and go to ι_{i+1} ; • $\tau !_k (k \neq j)$: if the pointed cell of tape τ is 0, go to ι_k ; else go to ι_{j+1} ; • τ_k^{k} $(k \neq j)$: if the pointed cell of tape τ is 1, go to ι_k ; else go to ι_{j+1} . ▶ Unary encoding of go-to's: (below, $\sigma \in \{\tau !, \tau ?\}_{\tau \in \{A, B\}}$; · is concatenation) • If k < j, encode σ_k as $\sigma \cdot -j^{-k} \cdot Q$ (i.e., repeating token "-" j - k times). • If k > j, encode σ_k as $\sigma \cdot +^{k-j} \cdot @$ (i.e., repeating token "+" k - j times).

TURING COMPLETENESS

a finite alphabet Σ , a <u>finite-size</u> decoder-only Transformer $\Gamma: \Sigma^+ \to \Sigma$, and coding schemes tokenize: $\{0,1\}^* \rightarrow \Sigma^*$ and readout: $\Sigma^* \rightarrow \{0,1\}^*$ with which prompting is **Turing-complete**, in the sense that • for every computable function φ : dom $\varphi \to \{0,1\}^*$ with dom $\varphi \subseteq \{0,1\}^*$, • there exists a prompt $\pi_{\varphi} \in \Sigma^+$ such that for every input $x \in \text{dom } \varphi$,

$$\mathbf{r}_{\varphi} \cdot \operatorname{tokenize}(\mathbf{x}) \Big) = \varphi(\mathbf{x}).$$

tokenize & readout run in time $O(|x|) \& O(|\varphi(x)|)$ on a RAM, respectively.

SIMULATION via Cot Steps

> Input tokenization:

> The constructed alphabet:

 $\Sigma = \{\#, AL, BL, AR, BR, A0, B0, A1, B1, A!, B!, A?, B?, -, +, @, ^, $, /, =, :, 0, 1\}.$ The construction of the Transformer is also omitted. See our paper...

 \succ **Example**: Suppose φ decides the DYCK language. DYCK: balanced parenthesis sequences; "0" as "(", "1" as ")". **Prompt** π_{ω} for deciding DYCK: ^A?++++++++++++#@A@ALA@ALA?---@ARARA1ARBLB?++@A1#ARA?++++@ B1BRB!+++@BLB!++++@B0ARB!-------------@ALARARA?--@ A0ALA0ALA?---@ARARA1#\$ > The input 00 has Shannon's encoding 1010, tokenized as: tokenize(00) = ARARARARALALA1ALA1=-----@ \succ The generated **CoT steps** for computing $\varphi(00)$: =+++++++++++++@AR/B1BR=+++@B0AR=-------@ =+++++++++++++@AR/B1BR=+++@B0AR=-------@ /AOALAOAL = - - - @AOALAOAL = - - - @AOALAOAL/ARARA1ARBL = + + @: 0\$The final readout answer is :0\$.

More examples on GitHub!

UIUC ILLINOIS 🛄

> Key idea: Record the execution of 2-PTMs via CoT steps.

 \succ Go-to instruction $ι_j = σ_k$ ($σ \in {\tau!, τ?}_{τ \in {A, B}}$): • If the go-to condition is not met, put token / in the CoT. • Else, put = $\cdot -j^{-k} \cdot @$ if k < j or = $\cdot +j^{-k} \cdot @$ if k > j.

> Halting & outputting: When execution reaches instruction #, • Put the output between two special tokens : and \$ in the CoT, • So that readout can extract it easily.

• Key idea: "Write" the input to tape A via CoT steps.

Trick: Encode the input *x* via **Shannon's** encoding.

Many details omitted here... See the example below.

DEMONSTRATION

It correctly computes $\varphi(00) = 0$ (00 \notin DYCK).

SCAN TO LEARN MORE

