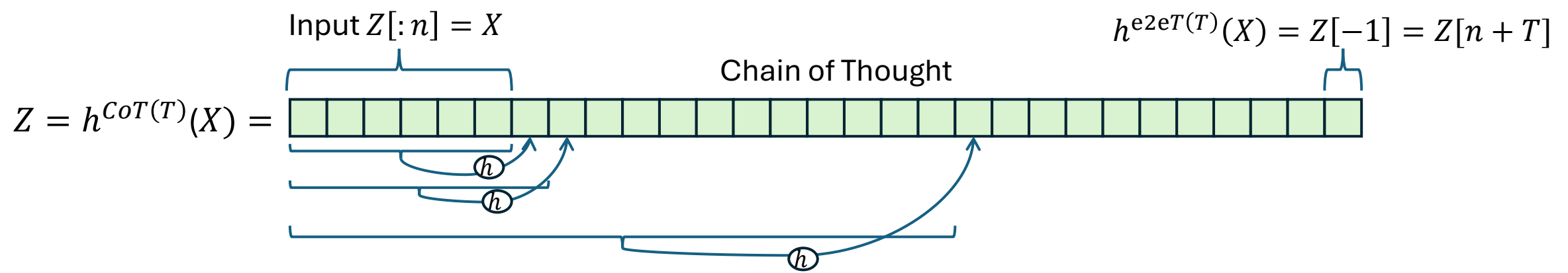


Learning Iterated Models

Nati Srebro (TTIC)

Based on ongoing joint work and discussions with
Nirmit Joshi (TTIC), **Gal Vardi** (TTIC→Weizmann),
Theo Misiakiewicz (TTIC→Yale), **Zhiyuan Li** (TTIC)



Iterating a function $h: \Sigma^* \rightarrow \Sigma$

Generally consider $h(X) = h([0] \oplus X)$, and think of $h(X) = h([0,0,0, \dots, 0] \oplus X)$ when defining $h(\cdot)$

For input of length $|X| = n$, start with $Z[0:n] = X$, and for $t \geq n$:

$$Z[t] = h(Z[:t])$$

$$h^{\text{CoT}(0)}(X) = X$$

$$h^{\text{CoT}(t+1)}(X) = h^{\text{CoT}(t)}(X) \oplus \left[h \left(h^{\text{CoT}(t)}(X) \right) \right] \quad \text{i.e. } h^{\text{CoT}(t)}(X) \in \Sigma^{|X|+t}$$

$$h^{\text{e2e}(t)}(X) = h^{\text{CoT}(t)}(X)[-1], \quad \text{i.e. the output of the last iteration}$$

For a base class $\mathcal{H} = \{h: \Sigma^* \rightarrow \Sigma\} \subseteq \Sigma^{\Sigma^*}$ consider:

$$\mathcal{H}^{\text{CoT}(T)} = \{ h^{\text{CoT}(T)}: \Sigma^* \rightarrow \Sigma^* \mid h \in \mathcal{H} \} \quad \mathcal{H}^{\text{e2e}(T)} = \{ h^{\text{e2e}(T)}: \Sigma^* \rightarrow \Sigma \mid h \in \mathcal{H} \} \subseteq \Sigma^{\Sigma^*}$$

(Realizable) Learning of an Iterated Class

End-to-end learning of $\mathcal{H}^{e2e(T)}$: Given $(X_i, y_i = h^{e2e(T)}(X_i))$ for $X_1 \dots X_m \sim iid \mathcal{D}$, learn $\hat{h}: \Sigma^* \rightarrow \Sigma$ with $\mathbb{E}_{X \sim \mathcal{D}} \left[err(\hat{h}(X), h^{e2e(T)}(X)) \right] \leq \epsilon$ $err(\hat{y}, y) = 1 \text{ iff } \hat{y} \neq y$

- Formally: Learning rule $A: (\Sigma^* \times \Sigma)^* \rightarrow \Sigma^*$ ϵ -e2e-learns $\mathcal{H}^{(T)}$ with sample complexity m on input distribution \mathcal{D} , if for any $h \in \mathcal{H}$, $\mathbb{E}_{S \sim \mathcal{D}_{h^{e2e(T)}}^m} \mathbb{E}_{X \sim \mathcal{D}} \left[err(A(S)(X), h^{e2e(T)}(X)) \right] \leq \epsilon$
- Dist independent learning: Rule A ϵ -e2e-learns with sample complexity $m(\epsilon)$ for any distribution \mathcal{D} over Σ^*
Input length dependent: sample complexity $m(\epsilon, n)$ for any distribution \mathcal{D} on Σ^n

\equiv PAC-Learning $\mathcal{H}^{e2e(T)}$

Chain-of-Thought learning of $\mathcal{H}^{e2e(T)}$: Given $\{Z_i = h^{CoT(T)}(X_i)\}$ for $X_1 \dots X_m \sim iid \mathcal{D}$, learn $\hat{h}: \Sigma^* \rightarrow \Sigma$ with $\mathbb{E}_{X \sim \mathcal{D}} \left[err(\hat{h}(X), h^{e2e(T)}(X)) \right] \leq \epsilon$

- Formally: Learning rule $A: (\Sigma^*)^* \rightarrow \Sigma^*$ ϵ -CoT-learns $\mathcal{H}^{(T)}$ with sample complexity m on input distribution \mathcal{D} , if for any $h \in \mathcal{H}$, $\mathbb{E}_{S \sim \mathcal{D}_{h^{CoT(T)}}^m} \mathbb{E}_{X \sim \mathcal{D}} \left[err(A(S)(X), h^{e2e(T)}(X)) \right] \leq \epsilon$

\approx AR-Learnable [Malach 24 “Auto-Regressive Next-Token Predictors are Universal Learners”]
(but we only care about final output)

Bounding the Sample Complexity in terms of \mathcal{H}

If \mathcal{H} has bounded cardinality:

$$m_{\text{CoT}(T)}(\mathcal{H}) \leq m_{e^{2e(T)}}(\mathcal{H}) \leq O\left(\log|\mathcal{H}^{e^{2e(T)}}| \cdot \frac{1}{\epsilon}\right) \leq O\left(\log|\mathcal{H}| \cdot \frac{1}{\epsilon}\right)$$

If $h \in \mathcal{H}$ has binary output, i.e. $\Sigma = \{0,1\}$, can we bound in terms of $VC(\mathcal{H})$?

$$m_{\text{CoT}(T)}(\mathcal{H}) \leq m_{e^{2e(T)}}(\mathcal{H}) \leq \tilde{O}\left(VC(\mathcal{H}^{e^{2e(T)}}) \cdot \frac{1}{\epsilon}\right) \leq \tilde{O}\left(T \cdot VC(\mathcal{H}) \cdot \frac{1}{\epsilon}\right)$$

Tight: $\forall_{D,T}$ exists \mathcal{H} with $VC(\mathcal{H}) = D$ but $m_{e^{2e(T)}}(\mathcal{H}) = \Omega\left(VC(\mathcal{H}^{e^{2e(T)}})\right) = \Omega(TD)$, even over $\{0,1\}^n$, $n = O(\log TD)$

Over $\Sigma = \mathbb{R}$ and relying on generalization of VC, even worse:

There exists \mathcal{H} with subgraph dimension (aka Pollard pseudo-dimension) 1, but $\mathcal{H}^{e^{2e(T)}}$ has infinite subgraph dim, for any $T > 1$.

Iterated Linear Thresholds

$$\mathcal{L}_d = \{ h_w: Z \mapsto \text{sign}(\langle w, [1] \oplus Z[-d:] \rangle) \mid w \in \mathbb{R}^{d+1} \} \text{ over } \Sigma = \{\pm 1\}$$

Theorem: $VC(\mathcal{L}_d^{e^{2e(T)}}) \leq O(dT \wedge d^2)$

e2e Learning

Find w s.t. $y_i = h_w^{e^{2e(T)}}(X_i)$
For $i = 1..m$

$$m = \tilde{O}\left(\left(dT \wedge d^2\right) \cdot \frac{1}{\epsilon}\right)$$

Open: is this tight?
Even $m = O(d)$ possible

Hard!

No $\text{poly}(n, d, T)$ time learning
(even improper, if \exists crypto)

CoT Learning

Find w s.t. $Z_i[-t] = h_w(Z_i[: -t])$
For $i = 1..m, t = 1..T$

With m contexts, mT total examples for \mathcal{H}

→ Suggests $m = \frac{d}{T} \cdot \frac{1}{\epsilon}$ might be enough

... but examples not iid (only d independent)

... AND need error at each step to be $\frac{\epsilon}{T}$

→ **Actual sample complexity** $\tilde{O}\left(\left(dT \wedge d^2\right) \cdot \frac{1}{\epsilon}\right)$

Easy!

LP-Satisfiability

Separate w_i (improper)

Find w_t s.t. $Z_i[-t] = h_{w_t}(Z_i[: -t])$
For $i = 1..m, t = 1..T$ [Malach 24]

m ind. Samples per w_i , but need $\text{err} \leq \frac{\epsilon}{T}$

→ **Sample complexity** $\tilde{O}\left(dT \cdot \frac{1}{\epsilon}\right)$

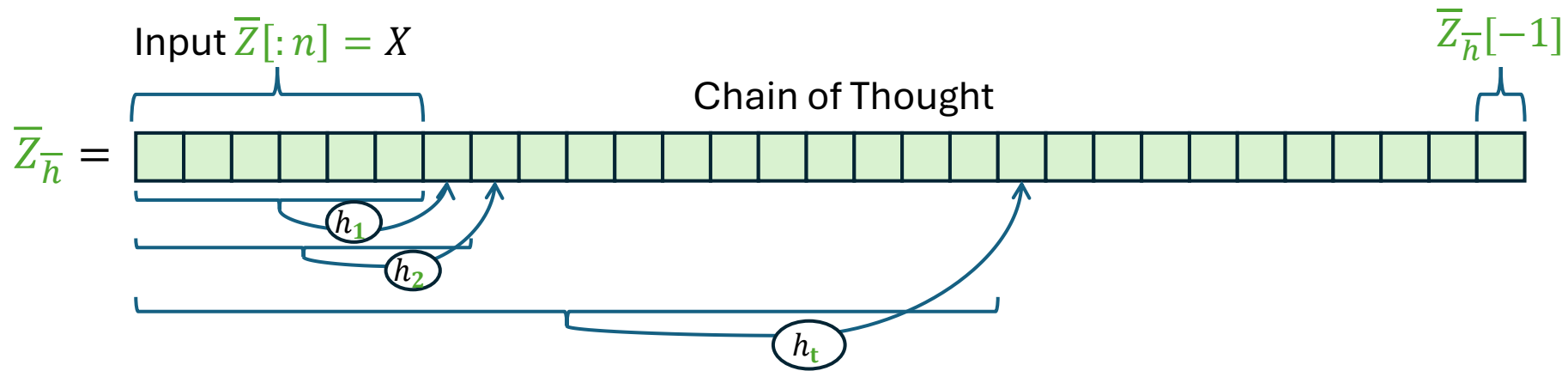
Open: is this tight?
No stat advantage over e2e?
Advantage over separate w_i ?

Simulating Circuits with Iterated Linear Thresholds

Recall $\mathbf{Z} = \mathbf{h}^{\text{CoT}(T)}(X)$ is defined as $\mathbf{Z}[t] = \mathbf{h}(X[:t])$

Contrast with $\bar{\mathbf{Z}}_{\bar{h}}$ specified by $\bar{h} = (h_1, \dots, h_T)$ defined as $\bar{\mathbf{Z}}[t] = h_t(\bar{\mathbf{Z}}[:t])$ initialized with $\bar{\mathbf{Z}}[:|X|] = X$

$$\mathcal{H}^{\text{CoT}(T)} = \{X \mapsto \bar{\mathbf{Z}}_{\bar{h}} \mid \bar{h} = (h_1, \dots, h_T) \in \mathcal{H}^T\}, \quad \mathcal{H}^{\text{e2e}(T)} = \{X \mapsto \bar{\mathbf{Z}}_{\bar{h}}[-1] \mid \bar{h} = (h_1, \dots, h_T) \in \mathcal{H}^T\}$$



Simulating Circuits with Iterated Linear Thresholds

Recall $\mathbf{Z} = \mathbf{h}^{\text{CoT}(T)}(X)$ is defined as $\mathbf{Z}[t] = \mathbf{h}(X[:t])$

Contrast with $\bar{\mathbf{Z}}_{\bar{\mathbf{h}}}$ specified by $\bar{\mathbf{h}} = (h_1, \dots, h_T)$ defined as $\bar{\mathbf{Z}}[t] = h_t(\bar{\mathbf{Z}}[:t])$ initialized with $\bar{\mathbf{Z}}[:|X|] = X$

$$\mathcal{H}^{\overline{\text{CoT}}(T)} = \left\{ X \mapsto \bar{\mathbf{Z}}_{\bar{\mathbf{h}}} \mid \bar{\mathbf{h}} = (h_1, \dots, h_T) \in \mathcal{H}^T \right\}, \quad \mathcal{H}^{\text{e2e}(T)} = \left\{ X \mapsto \bar{\mathbf{Z}}_{\bar{\mathbf{h}}}[-1] \mid \bar{\mathbf{h}} = (h_1, \dots, h_T) \in \mathcal{H}^T \right\}$$

Any (logical or linear threshold) circuit of size S over $\{\pm 1\}^n$ can be simulated by $\mathcal{L}_{n+T}^{\overline{\text{CoT}}(T=S)}$,
 i.e. for any circuit, there is a sequence $\bar{\mathbf{h}}$ of lin thresholds s.t. $\bar{\mathbf{Z}}_{\bar{\mathbf{h}}}[i]$ is the output of unit i on input $\bar{\mathbf{Z}}_{\bar{\mathbf{h}}}[:n] = X$
 [Malach 24]

Theorem: For any $w_1, \dots, w_T \in \mathbb{R}^d$, i.e. $\bar{\mathbf{h}} \in \mathcal{L}_d^T$, there exists $\mathbf{w} \in \mathbb{R}^{\mathcal{O}((n+T)^2)}$ s.t. for all $X \in \{\pm 1\}^n$,

$$\mathbf{h}_{\mathbf{w}}^{\text{e2e}(\mathcal{O}((n+T)^2))} \left(\begin{bmatrix} -1 & 1^{\oplus(n+T)^2} \end{bmatrix} \oplus X \right) = \bar{\mathbf{Z}}_{\bar{\mathbf{h}}}(X)[-1]$$

Conclusion: Any circuit of size S can be simulated by $\mathcal{L}_d^{\text{e2e}(T)}$ with $d, T = \mathcal{O}(S^2)$ and a fixed input expansion.

Learning poly-size circuits is hard \rightarrow learning $\mathcal{L}_d^{\text{e2e}(T)}$ in time $\text{poly}(n, d, T)$ is hard

What I'd Really Like to Learn

- Ultimate: $DESC_S = \{ \text{any function describable with } S \text{ bits} \}$
 - $\log|DESC_S| \leq S$, hence learnable with $O(S)$ samples
- More specifically: $PROG_S = \{ \text{programs with of length } S \}$
 - $\log|PROG_S| = O(S)$, hence learnable with $O(S)$ samples
 - But: Not tractable computationally + output not useful
- How about: $TIME_T = \{ \text{programs with runtime } T \}$
 - Also learnable with $O(T)$ samples, learning $\in NP$
 - $TIME_T \subseteq CIRCUT_{poly(T)} \subseteq NN_{poly(T)}$
 - We can learn $TIME_T$ (with poly-opt sample size) by learning a Feed-Forward Neural Nets
- But what about $PROG_{S,T} = \{ \text{programs of length } S \text{ and runtime } T \}$
 - Learning with $O(S)$ samples is $\in NP$, as long as $T = poly(n)$
 - “Goal”: learn with $poly(S \log T)$ samples and “training time” $poly(T)$
 - Can we construct a simple and “trainable” class $\mathcal{H} \supseteq PROG_{S,T}$ of complexity $poly(S \log T)$?

$$\mathcal{L}_d = \{h_w: Z \mapsto \text{sign}(\langle w, [1] \oplus Z[-d:] \rangle) \mid w \in \mathbb{R}^{d+1}\}$$

$$Z = h_w^{\text{CoT}}(Z[:n]) \text{ with } Z[t] = h_w(Z[:t]) \text{ for } t \geq n$$

Problem with $\mathcal{L}_d^{\text{CoT}(T)}$: context length is $O(d)$

→ i.e. $Z[t]$ only depends on $Z[t-d:t] \in \{\pm 1\}^d$

→ Computation has state space of size 2^d , i.e. need $d = \Omega(\text{MEMORY})$

How can we get context length \gg complexity?

Simplest attempt: sparse linear thresholds:

$$\mathcal{L}_{d,k} = \{h_w: Z \mapsto \text{sign}(\langle w, [1] \oplus Z[-d:] \rangle) \mid w \in \mathbb{R}^{d+1}, \|w\|_0 \leq k\}$$

- Complexity $VC\left(\mathcal{L}_{d,k}^{e^{2e(T)}}\right) = O(k^2 \log d)$ with context length d
- Tractably CoT learnable for $k = O(1)$ (is this enough?) or with a sparse-learning oracle.
- What is the expressive power of $\mathcal{L}_{d,k}^{\text{CoT}(T)}$??

- Can we describe a simple class \mathcal{H} s.t.:
 - $PROG_{S,T} \subseteq \mathcal{H}^{e^{2e(T)}}$
 - $VC(\mathcal{H}^{e^{2e(T)}}) = poly(S \log T)$
 - CoT learnable with $poly(S \log T)$ and time $poly(T)$

The Turing Iterative Model

For an internal state space size S (and wlog tape alphabet $\{0,1\}$) consider a iterative model on alphabet $\Sigma = [S] \times \{0,1\} \times \{-1,0,1\}$

parametrized by a Turing Machine transition function $A: (state, tape\ sym) \mapsto (new\ state, new\ sym, move)$:
 $\mathcal{H}_S = \{ h_A: \Sigma^* \rightarrow \Sigma \mid A \in ([S] \times \{0,1\} \times \{-1,0,1\})^{[S] \times \{0,1\}} \}$

Where $h_A((s_0, r_0, m_0), \dots, (s_{t-1}, r_{t-1}, m_{t-1}))$ computes the next symbol to be written, state and move for a TM specified by transition table A

$h_A(Z = (s_0, r_0, m_0), \dots, (s_{t-1}, r_{t-1}, m_{t-1}))$:
For each $i = 0 \dots t$, calculate $pos[i] = \sum_{j < i} S_j$
Find $i^* = \max i$ s.t. $pos[i] = pos[t]$
If exists, $r = r_{i^*}$, else $r = 0$
output $(s_t, r_t, m_t) = A(s_{t-1}, r)$

Claim: $PROG_{poly(S), poly(T)} = TM_{S,T} \subseteq \mathcal{H}_S^{e2e(T)}$
With a fixed input embedding $Z[i] = (0, X[i], 1)$
And fixed output decoding (proj to 2nd component)

Claim: $\log |\mathcal{H}_S| \leq 2S \log 6S$
 \rightarrow e2e learnable with $m = O(S \log S)$

Claim: CoT learnable with $m = O(S \log S)$ in time $\tilde{O}(ST)$
by memorizing size- $O(S)$ table

The Turing Iterative Model

For an internal state space size S (and wlog tape alphabet $\{0,1\}$) consider an iterative model on alphabet $\Sigma = [S] \times \{0,1\} \times \{-1,0,1\}$

parametrized by a Turing Machine transition function $A: (state, tape\ sym) \mapsto (new\ state, new\ sym, move)$:
 $\mathcal{H}_S = \{ h_A: \Sigma^* \rightarrow \Sigma \mid A \in ([S] \times \{0,1\} \times \{-1,0,1\})^{[S] \times \{0,1\}} \}$

Where $h_A((s_0, r_0, m_0), \dots, (s_{t-1}, r_{t-1}, m_{t-1}))$ computes the next symbol to be written, state and move for a TM specified by transition table A

$h_A(Z = (s_0, r_0, m_0), \dots, (s_{t-1}, r_{t-1}, m_{t-1}))$:

For each $i = 0 \dots t$, calculate $pos[i] = \sum_{j < i} s_j$ $pos[i] = i \cdot \text{Attn}(k[j] = 0, q[i] = 0, v[j] = s_j)$

Find $i^* = \max i$ s.t. $pos[i] = pos[t]$

If exists, $r = r_{i^*}$, else $r = 0$

output $(s_t, r_t, m_t) = A(s_{t-1}, r)$

Keys: $k[i] = (i, 100T \ pos[i])$
 Values: $v[i] = (pos[i], r_i)$
 Query: $q[t] = (t, 100T \ pos[t])$
 Obtain $(p, r) = \text{Attn}(k, q[t], v)$
 If $p \neq pos[i]$, set $r \leftarrow 0$

$\text{Attn}(k, q[i], v)$:
 $j^* = \arg \min_{j < i} \|k[j] - q[i]\|$
 return $v[j^*]$

The Turing Iterative Model

For an internal state space size S (and wlog tape alphabet $\{0,1\}$) consider an iterative model on alphabet $\Sigma = [S] \times \{0,1\} \times \{-1,0,1\}$

parametrized by a Turing Machine transition function $A: (state, tape\ sym) \mapsto (new\ state, new\ sym, move)$:
 $\mathcal{H}_S = \{ h_A: \Sigma^* \rightarrow \Sigma \mid A \in ([S] \times \{0,1\} \times \{-1,0,1\})^{[S] \times \{0,1\}} \}$

Where $h_A((s_0, r_0, m_0), \dots, (s_{t-1}, r_{t-1}, m_{t-1}))$ computes the next symbol to be written, state and move for a TM specified by transition table A

$h_A(Z = (s_0, r_0, m_0), \dots, (s_{t-1}, r_{t-1}, m_{t-1}))$:

For each $i = 0 \dots t$, calculate $pos[i] = \sum_{j < i} s_j$

$pos[i] = i \cdot \text{Attn}(k[j] = 0, q[i] = 0, v[j] = s_j)$

Find $i^* = \max i$ s.t. $pos[i] = pos[t]$

If exists, $r = r_{i^*}$, else $r = 0$

output $(s_t, r_t, m_t) = A(s_{t-1}, r)$

Keys: $k[i] = (i, 100T \cdot pos[i])$

Values: $v[i] = (pos[i], r_i)$

Query: $q[t] = (t, 100T \cdot pos[t])$

Obtain $(p, r) = \text{Attn}(k, q[t], v)$

If $p \neq pos[i]$, set $r \leftarrow 0$

$\text{Attn}(k, q[i], v)$:

$w[j] = 1$ iff

$j \in \arg \min_{j < i} \|k[j] - q[i]\|$

return $\frac{\sum_{j < i} w[j] v[j]}{\sum_{j < i} w[j]}$

The Turing Iterative Model

For an internal state space size S (and wlog tape alphabet $\{0,1\}$) consider an iterative model on alphabet $\Sigma = [S] \times \{0,1\} \times \{-1,0,1\}$

parametrized by a Turing Machine transition function $A: (state, tape\ sym) \mapsto (new\ state, new\ sym, move)$:
 $\mathcal{H}_S = \{ h_A: \Sigma^* \rightarrow \Sigma \mid A \in ([S] \times \{0,1\} \times \{-1,0,1\})^{[S] \times \{0,1\}} \}$

Where $h_A((s_0, r_0, m_0), \dots, (s_{t-1}, r_{t-1}, m_{t-1}))$ computes the next symbol to be written, state and move for a TM specified by transition table A

$h_A(Z = (s_0, r_0, m_0), \dots, (s_{t-1}, r_{t-1}, m_{t-1}))$:

For each $i = 0 \dots t$, calculate $pos[i] = \sum_{j < i} s_j$ $\leftarrow pos[i] = i \cdot \text{Attn}(k[j] = 0, q[i] = 0, v[j] = s_j)$

Find $i^* = \max i$ s.t. $pos[i] = pos[t]$

If exists, $r = r_{i^*}$, else $r = 0$

output $(s_t, r_t, m_t) = A(s_{t-1}, r)$

Keys: $k[i] = (i, 100T \cdot pos[i])$
 Values: $v[i] = (pos[i], r_i)$
 Query: $q[t] = (t, 100T \cdot pos[t])$
 Obtain $(p, r) = \text{Attn}_{\beta \rightarrow \infty}(k, q[t], v)$
 If $p \neq pos[t]$, set $r \leftarrow 0$

$\text{Attn}_{\beta}(k, q[i], v)$:
 $w[j] = e^{\beta \|k[j] - q[i]\|}$
 return $\frac{\sum_{j < i} w[j] v[j]}{\sum_{j < i} w[j]}$

The Turing Iterative Model

For an internal state space size S (and wlog tape alphabet $\{0,1\}$) consider an iterative model on alphabet $\Sigma = [S] \times \{0,1\} \times \{-1,0,1\}$

parametrized by a Turing Machine transition function $A: (state, tape\ sym) \mapsto (new\ state, new\ sym, move)$:
 $\mathcal{H}_S = \{ h_A: \Sigma^* \rightarrow \Sigma \mid A \in ([S] \times \{0,1\} \times \{-1,0,1\})^{[S] \times \{0,1\}} \}$

Where $h_A((s_0, r_0, m_0), \dots, (s_{t-1}, r_{t-1}, m_{t-1}))$ computes the next symbol to be written, state and move for a TM specified by transition table A

$h_A(Z = (s_0, r_0, m_0), \dots, (s_{t-1}, r_{t-1}, m_{t-1}))$:

For each $i = 0 \dots t$, calculate $pos[i] = \sum_{j < i} s_j$ $\leftarrow pos[i] = i \cdot \text{Attn}(k[j] = 0, q[i] = 0, v[j] = s_j)$

Find $i^* = \max i$ s.t. $pos[i] = pos[t]$

If exists, $r = r_{i^*}$, else $r = 0$

output $(s_t, r_t, m_t) = A(s_{t-1}, r)$

Keys: $k[i] = (i, pos[i]^2, T^2 - pos[i]^2)$

Values: $v[i] = (pos[i], r_i)$

Query: $q[t] = \left(\frac{1}{100T^4}, pos[t]^2, T^2 - pos[t]^2 \right)$

Obtain $(p, r) = \text{Attn}_{\beta \rightarrow \infty}(k, q[t], v)$

If $p \neq pos[i]$, set $r \leftarrow 0$

$\text{Attn}(k, q[i], v)$:

$w[j] = e^{\beta \langle q[i], v \rangle}$

return $\frac{\sum_{j < i} w[j] v[j]}{\sum_{j < i} w[j]}$

$\mathcal{H}_S \subseteq \{ \text{Transformer with } \text{poly}(S \log T)\text{-sized MLP and } O(\log T)\text{ precision} \}$ (with simple fixed encoding)

[Jorge Pérez, Pablo Barceló, Javier Marinkovic '21, Attention is Turing-Complete]

[Colin Wei, Yining Chen, Tengyu Ma '22, Statistically Meaningful...Approximating TM with Transformers]

[William Merrill, Ashish Sabharwal '24, The Expressive Power of Transformers with Chain of Thought]

A (Minimal) Transformer?

$$f_w(X) = \begin{cases} \text{sign}(\langle w, Z[-d:] \rangle) & \text{if } Z[-5:] \neq [1 \ 1 \ 1 \ 1 \ 1] \\ \text{Attn}(Z) & \text{if } Z[-5:] = [1 \ 1 \ 1 \ 1 \ 1] \end{cases}$$

$$\mathcal{H}_{d,r} = \{ f_w \mid w \in \mathbb{R}^d \}$$

Attn(X):

$$q = X[-r:0]$$

$$k[i] = X[-2ri: -2ri + r]$$

$$v[i] = X[-2ri + r]$$

return **Attn**(keys = k, query = q, vals = v)

Claim: $\mathcal{H}_{d,r}^{(T)}$ CoT learnable with $m = O(d^2)$ samples in time $\text{poly}(d, T)$

What can it represent? $\text{PROG}_{S,T} \in \mathcal{H}_{d,r=\text{poly}(S \log T)}^{\text{e2e}(\text{poly}(T))}$???

Summary

- Study of (Stationary) Iterative Models essential for understanding:
 - Autoregressive learning as it is actually done
 - Learning with large context length, and sample complexity independent of context (input + output) length
 - Length generalization
 - Learning with sample complexity scaling with program length, not runtime
- Open questions even on simple models
- What's the right view of non-realizable learning?
 - What is the goal/reward/error?
 - Discriminative/reward-based view of iterative models vs generative view