# Comp-mech for logits:

## Modular addition as a case study

Xavier Poncini

PIBBSS Symposium 2025

# Outline:

1) (Just enough) Comp-mech

2) Modular addition

3) Results

4) Outlook

(Just enough) <u>Comp–Mech</u>

# Hidden Markov Model (HMM):

A HMM consists of:

* A set $\mathcal{X}$ ~~~ "think" ~~~ vocabulary of emissions

* A collection of transition matrices $(T^{(x)})_{x \in \mathcal{X}}$ ~~~ "think" ~~~ dynamic that determines emissions & hidden state transitions

where $T_{ij}^{(x)} = P(X=x, \; S_j = s_j \mid S_i = s_i)$

Fix an initial dist. over hidden states $(P(s_i))$; the prob. of obs. $W = w_1 w_2 \ldots w_n$ is:

$$P(W) = \sum_{i,j,\ldots,\ell} P(s_i) \, P(w_1, s_j \mid s_i) \, P(w_2, s_k \mid s_j) \cdots P(w_n, s_\ell \mid s_j)$$

$$= \langle \eta | \; T^{(w_1)} T^{(w_2)} \ldots T^{(w_n)} \; | \tau \rangle \qquad \text{where} \quad \langle \eta | = [\, P(s_1) \; \ldots \; P(s_{|S|}) \,]$$

$$\underbrace{\phantom{T^{(w_1)} T^{(w_2)} \ldots T^{(w_n)}}}_{T^{(W)}} \qquad | \tau \rangle = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$\underline{Example:}$$

$$\mathcal{X} = \{0, 1\}$$

$$T^{(0)} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ \frac{1}{2} & 0 & 0 \end{bmatrix}$$

$$T^{(1)} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ \frac{1}{2} & 0 & 0 \end{bmatrix}$$

What information about the _past_ is relevant to predict the _future_?

Consider conditional probabilities:

$$P(w^{(f)} \mid w^{(p)}) = \frac{P(w^{(p)} w^{(f)})}{P(w^{(p)})} = \frac{\langle \eta | T^{(w^{(p)})} T^{(w^{(f)})} | 1 \rangle}{\langle \eta | T^{(w^{(p)})} | 1 \rangle}$$

The _predictive vector_ $\langle \eta^{(w^{(p)})} | = \dfrac{\langle \eta | T^{(w^{(p)})}}{\langle \eta | T^{(w^{(p)})} | 1 \rangle}$ is relevant to _all_ future predictions:

$$P(w^{(f)} \mid w^{(p)}) = \langle \eta^{(w^{(p)})} | T^{(w^{(f)})} | 1 \rangle$$

## Observation 1: [Shai et al.]

The **predictive vector** of the HMM is <u>often</u> linearly decodable from the activations of a neural network trained on data from the HMM.

## Observation 2:

Transformers (& other NNs) produce probabilities by passing logits through an "<u>intense</u>" non-linearity — the softmax function:

(Boltzmann dist.)

$$P(w) = \frac{e^{z(w)}}{\sum_w e^{z(w)}}$$

## Question:

Can we develop a notion of a HMM for logits that admits an analogue of predictive vectors?

# Energy-based hidden Markov model (EHMM):

A EHMM consists of:

* A set $\mathcal{X}$

* A collection of transition matrices $(H^{(x)})_{x \in \mathcal{X}}$

* An initial vector $\langle \eta |$ & a final vector $| \varphi \rangle$

Such that $\langle \eta | H^{(w)} | \varphi \rangle \in \mathbb{R}$, for all $w \in \mathcal{X}^N$ & $N \in \mathbb{N}$

$$H^{(w_1)} H^{(w_2)} \ldots H^{(w_N)}$$

We can then associate these matrix elements with logits (energies)

$$Z(w) = \langle \eta | H^{(w)} | \varphi \rangle \quad \leadsto \quad P(w) = \frac{e^{z(w)}}{\sum_w e^{z(w)}}$$

What information about the past is relevant to predict the future?

Consider "conditional logits" i.e. $Z(W^{(f)} | W^{(p)})$ such that:

$$P(W^{(f)} | W^{(p)}) = \frac{e^{Z(W^{(f)} | W^{(p)})}}{\sum_{W^{(f)}} e^{Z(W^{(f)} | W^{(p)})}}$$

Claim: $Z(W^{(f)} | W^{(p)}) = Z(W^{(p)} W^{(f)})$   (proof is easy)

Expressing conditional logits in terms of EHMM objects:

$$Z(W^{(f)} | W^{(p)}) = Z(W^{(p)} W^{(f)}) = \langle \eta | H^{(W^{(p)})} H^{(W^{(f)})} | e \rangle$$

The predictive vector $\langle \eta^{(W^{(p)})} | = \langle \eta | H^{(W^{(p)})}$ is relevant to all future predictions:

$$Z(W^{(f)} | W^{(p)}) = \langle \eta^{(W^{(p)})} | H^{(W^{(f)})} | 1 \rangle$$

## Summary:

| Process | Output | Predictive vector |
|---------|--------|-------------------|
| HMM | $P(w) = \langle \eta \mid T^{(w)} \mid 1 \rangle$ | $\langle \eta^{(w^{(p)})} \mid = \dfrac{\langle \eta \mid T^{(w^{(p)})}}{\langle \eta \mid T^{(w^{(p)})} \mid 1 \rangle}$ |
| EHMM | $Z(w) = \langle \eta \mid H^{(w)} \mid 1 \rangle$ | $\langle \eta^{(w^{(p)})} \mid = \langle \eta \mid H^{(w^{(p)})}$ |

## Questions:

* Do neural networks represent the predictive vector of the EHMM?

* Do neural networks prefer the predictive vector of the HMM over the EHMM when given the chance?

# Modular addition

$$a + b = c \mod p$$

# Cyclic group ($C_p$):

The group $C_p$ is generated by $r$ subject to:

$$r^p = id$$

E.g. $r^{a+b} = r^c \iff c = a+b \bmod p$.

A $C_p$-action describes the action of $C_p$ on a set:

$$\alpha : S \times C_p \longrightarrow S$$

That respects the group structure i.e.

$$\alpha(v, r^{a+b}) = \alpha(\alpha(v, r^a), r^b).$$

Consider two $C_p$-actions: [Chughtai et al.]  $[\underset{0^{th}}{0} \cdots 0 \underset{i^{th}}{1} 0 \cdots \underset{(p-1)^{th}}{0}]$

1) $S_p \times C_p \longrightarrow S_p$, $\quad S_p = \{\langle e_i | \, | \, i = 0, 1, \ldots, p-1\}$

on the vertices of a $(p-1)$-simplex $S_p$. Inducing:

$$e : C_p \longrightarrow Mat_{p \times p}(\{0,1\}), \quad e^{(r)} = \begin{bmatrix} 0 & 1 & 0 \cdots & 0 \\ \vdots & & 0 & \ddots & 0 \\ 0 & & & 0 & 1 \\ 1 & 0 \cdots & & 0 \end{bmatrix}$$

$[\cos(\frac{2\pi \omega i}{p}) \quad \sin(\frac{2\pi \omega i}{p})]$

2) $V_p^{(\omega)} \times C_p \longrightarrow V_p^{(\omega)}$, $\quad V_p^{(\omega)} = \{\langle v_i^{(\omega)} | \, | \, i = 0, 1, \ldots, p-1\}$

on the vertices of a $p$-gon $V_p$. Inducing:

$$e_\omega : C_p \longrightarrow Mat_{2 \times 2}(\mathbb{R}), \quad e^{(r)} = \begin{bmatrix} \cos(\frac{2\pi \omega}{p}) & \sin(\frac{2\pi \omega}{p}) \\ -\sin(\frac{2\pi \omega}{p}) & \cos(\frac{2\pi \omega}{p}) \end{bmatrix}$$

# Random — Random Mod $p$ (RRModp):

Vocabulary: $\mathcal{X} = \{0, 1, \ldots, p-1\}$

Hidden states: $S = \{s_0^{(0)}\} \cup \{s_0^{(1)}, \ldots, s_{p-1}^{(1)}\} \cup \{s_0^{(2)}, \ldots, s_{p-1}^{(2)}\}$

Intuitively, the RRModp HMM is given by:

0) Process initialised in state $s_0^{(0)}$

1) Sample $a$ from $\mathcal{X}$ & transition $s_0^{(0)} \xrightarrow{a} s_a^{(1)}$

2) Sample $b$ from $\mathcal{X}$ & transition $s_a^{(1)} \xrightarrow{b} s_{a+b \bmod p}^{(2)}$
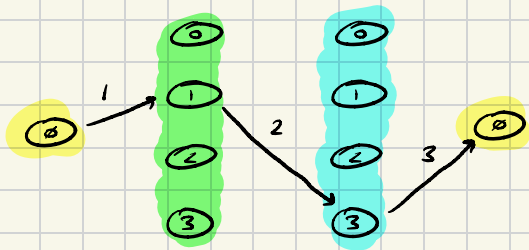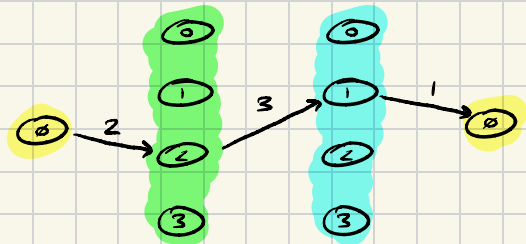
3) Transition $s_{a+b \bmod p}^{(2)} \xrightarrow{c} s_0^{(0)}$  where  $c = a+b \bmod p$

# Examples: $p = 4$



"a"     "b"     "c"

Formally, the RRMod$p$ HMM is defined:

$$\mathcal{K} = \{0, 1, \ldots, p-1\}, \quad T^{(i)} = \begin{array}{ccc} S^{(0)} & S^{(1)} & S^{(2)} \end{array}$$

$$\langle \mathcal{M}| = \begin{bmatrix} 1 & \underset{\sim}{0} & \underset{\sim}{0} \end{bmatrix}$$

$$T^{(i)} = \begin{array}{c} S^{(0)} \\ S^{(1)} \\ S^{(2)} \end{array} \begin{bmatrix} 0 & \frac{1}{p}\langle e_i | & \underset{\sim}{0} \\ \underset{\sim}{0} & \underset{\sim}{0} & \frac{1}{p} e(v^i) \\ |e_i\rangle & \underset{\sim}{0} & \underset{\sim}{0} \end{bmatrix}$$

Probabilities:

$$P(a) = \langle \mathcal{M}| T^{(a)} |\mathcal{1}\rangle = \begin{bmatrix} 0 & \frac{1}{p}\langle e_a| & \underset{\sim}{0} \end{bmatrix} \begin{bmatrix} \vdots \\ 1 \\ 1 \end{bmatrix} = \frac{1}{p} \qquad \langle e_a | e(r^b)$$

$$P(ab) = \langle \mathcal{M}| T^{(a)} T^{(b)} |\mathcal{1}\rangle = \begin{bmatrix} 0 & \underset{\sim}{0} & \frac{1}{p^2}\langle e_{a+b \bmod p}| \end{bmatrix} \begin{bmatrix} \vdots \\ 1 \\ 1 \end{bmatrix} = \frac{1}{p^2}$$

$$P(abc) = \langle \mathcal{M}| T^{(a)} T^{(b)} T^{(c)} |\mathcal{1}\rangle = \begin{bmatrix} \frac{1}{p^2}\langle e_{a+b \bmod p}| e_c\rangle & \underset{\sim}{0} & \underset{\sim}{0} \end{bmatrix} \begin{bmatrix} \vdots \\ 1 \\ 1 \end{bmatrix} = \frac{\delta_{a+b \bmod p, c}}{p^2}$$

## Predictive vectors of RR Mod p:

Recall the HMM predictive vector: $\langle \eta^{(w)} | = \dfrac{\langle \eta | T^{(w)}}{\langle \eta | T^{(w)} | \uparrow \rangle}$

Reusing parts of previous calculations:

$$\langle \eta^{(a)} | = \frac{\langle \eta | T^{(a)}}{\langle \eta | T^{(a)} | \uparrow \rangle} = [\, 0 \;\; \langle e_a | \;\; 0 \,]$$

$$\langle \eta^{(ab)} | = \frac{\langle \eta | T^{(a)} T^{(b)}}{\langle \eta | T^{(a)} T^{(b)} | \uparrow \rangle} = [\, 0 \;\; 0 \;\; \langle e_{a+b \bmod p} | \,]$$

$$\langle \eta^{(abc)} | = \frac{\langle \eta | T^{(a)} T^{(b)} T^{(c)}}{\langle \eta | T^{(a)} T^{(b)} T^{(c)} | \uparrow \rangle} = \begin{cases} \langle \eta | & , \; c = a+b \bmod p \\ \text{undefined} & , \; \text{else} \end{cases}$$

Projecting out the zero entries, the set of predictive vectors is given by:

$$S_p = \left\{ \langle e_i | \;\middle|\; i = 0, 1, \ldots, p-1 \right\}$$

$\underbrace{[\, 0 \cdots 0 \;\; 1 \;\; 0 \cdots 0 \,]}_{}$
$0^{th}$    $i^{th}$    $(p-1)^{th}$

vertices of a
$(p-1)$-simplex

## Soft Random - Random Mod p (sRRModp):

Vocabulary: $\mathcal{X} = \{0, 1, \ldots, p-1\}$

Vectors: $\langle v_i^{(\omega)} | = \left[ \cos\left(\frac{2\pi \omega i}{p}\right) \quad \sin\left(\frac{2\pi \omega i}{p}\right) \right]$

Intuitively, the sRRModp EHMM is given by:

0) Process initialised in vector $\langle q | = [1 \quad 0 \quad 0]$

1) Sample $a$ from $\mathcal{X}$ & transition $\langle q | \xrightarrow{a} \langle v_a^{(\omega)} |$

2) Sample $b$ from $\mathcal{X}$ & transition $\langle v_a^{(\omega)} | \xrightarrow{b} \langle v_{a+b}^{(\omega)} |$

3) Transition $\langle v_{a+b}^{(\omega)} | \xrightarrow{c} \langle q | \quad \text{argmax } P(* \, | \, a b) = c$

$\qquad\qquad$ where $c = a + b \mod p$

Fix a tuple of frequencies $\underset{\sim}{\omega} = (\omega_1, \omega_2, \ldots, \omega_N)$ where
$\omega_i \in \left\{ 1, 2, \ldots, \lfloor \frac{p}{2} \rfloor \right\}$.

Formally, the sRRMod$p$ EHMM is defined:

$$\mathcal{K} = \{0, 1, \ldots, p-1\},$$

$$H^{(i)} = \begin{bmatrix} 0 & \frac{1}{p} \overset{N}{\underset{j=1}{\bigoplus}} \langle v_i^{(\omega_i)} | & \underset{\sim}{0} \\ \underset{\sim}{0} & \underset{\sim}{0} & \frac{1}{p} \overset{N}{\underset{j=1}{\bigoplus}} e_{\omega_j}(r) \\ \overset{N}{\underset{j=1}{\bigoplus}} |v_i^{(\omega_i)}\rangle & \underset{\sim}{0} & \underset{\sim}{0} \end{bmatrix}$$

$$\langle \eta | = \begin{bmatrix} 1 & \underset{\sim}{0} & \underset{\sim}{0} \end{bmatrix}, \quad |\varphi\rangle = \begin{bmatrix} 1 \\ \underset{\sim}{0} \\ \underset{\sim}{0} \end{bmatrix}$$

## Logits:

$$z(a) = \langle \eta | H^{(a)} | \varphi \rangle = \left[ 0 \quad \frac{1}{p} \bigoplus_{j=1}^{N} \langle v_a^{(w_j)} | \quad \underset{\sim}{0} \right] \begin{bmatrix} 1 \\ 0 \\ \underset{\sim}{0} \\ \underset{\sim}{0} \end{bmatrix} = 0$$

$$z(ab) = \langle \eta | H^{(a)} H^{(b)} | \varphi \rangle = \left[ 0 \quad \underset{\sim}{0} \quad \frac{1}{p^2} \bigoplus_{j=1}^{N} \langle v_{a+b \bmod p}^{(w_j)} | \right] \begin{bmatrix} 1 \\ 0 \\ \underset{\sim}{0} \\ \underset{\sim}{0} \end{bmatrix} = 0$$

$$z(abc) = \langle \eta | H^{(a)} H^{(b)} H^{(c)} | \varphi \rangle$$

$$= \left[ \frac{1}{p^2} \sum_{j=1}^{N} \langle v_{a+b \bmod p}^{(w_j)} | v_c^{(w_j)} \rangle \quad \underset{\sim}{0} \quad \underset{\sim}{0} \right] \begin{bmatrix} 1 \\ 0 \\ \underset{\sim}{0} \\ \underset{\sim}{0} \end{bmatrix} = \frac{1}{p^2} \sum_{j=1}^{N} \cos\left( \frac{2\pi w_j (a+b-c)}{p} \right)$$

[Nanda et al.]
logits

When $c = a+b \bmod p$ cosines <u>constructively</u>
<u>interfere</u> & $z(abc)$ is <u>large</u>.

When $c \neq a+b \bmod p$ cosines <u>destructively</u>
<u>interfere</u> & $z(abc)$ is <u>small</u>.

## Probabilities:

We have 
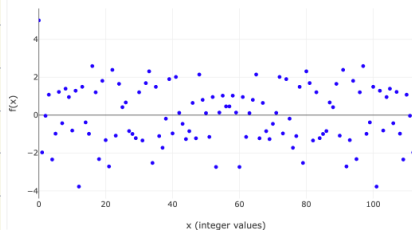$$P(w) = \frac{e^{z(w)}}{\sum\limits_{w \in \chi^L} e^{z(w)}} \quad \text{and}:$$

$$z(a) = z(ab) = 0 \;, \quad z(abc) = \frac{1}{p^2} \sum_{j=1}^{N} \cos\left(\frac{2\pi w_j \, (a+b-c)}{p}\right)$$

So 
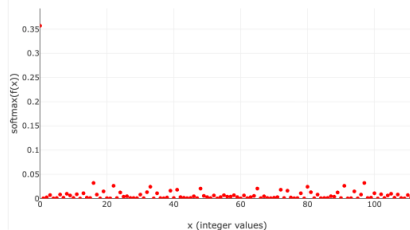$$P(a) = \frac{1}{p} \;, \quad P(ab) = \frac{1}{p^2} \quad \text{and}$$

$$P(abc) = \frac{\exp\left(\frac{1}{p^2} \sum\limits_{j=1}^{N} \cos\left(\frac{2\pi w_j \, (a+b-c)}{p}\right)\right)}{\sum\limits_{a,b,c} \exp\left(\frac{1}{p^2} \sum\limits_{j=1}^{N} \cos\left(\frac{2\pi w_j \, (a+b-c)}{p}\right)\right)}$$

# Interactive Discrete Cosine Sum & Softmax Visualization

### Discrete Function: $\Sigma a_i \cos(2\pi\omega_i x/p)$



### Discrete Softmax of Function



**N (terms):** 5 ⌄  **p:** 113

**Term 1**
α1: ———●——— 1   ω1: ——●———— 14

**Term 2**
α2: ———●——— 1   ω2: ————●—— 35

**Term 3**
α3: ———●——— 1   ω3: —————●— 41

**Term 4**
α4: ——●———— 1   ω4: —————●— 42

**Term 5**
α5: ———●——— 1   ω5: —————●— 52

## Predictive vectors of sRR Mod p:

Recall the EHMM predictive vector: $\langle \eta^{(w)}| = \langle \eta | H^{(w)}$

Reusing parts of previous calculations:

$$\langle \eta^{(a)}| = \langle \eta | H^{(a)} = \left[ 0 \quad \frac{1}{p} \bigoplus_{j=1}^{N} \langle v_a^{(w_j)}| \quad \underset{\sim}{0} \right]$$

$$\langle \eta^{(ab)}| = \langle \eta | H^{(a)} H^{(b)} = \left[ 0 \quad \underset{\sim}{0} \quad \frac{1}{p^2} \bigoplus_{j=1}^{N} \langle v_{a+b \bmod p}^{(w_j)}| \right]$$

$$\langle \eta^{(abc)}| = \langle \eta | H^{(a)} H^{(b)} H^{(c)} = \left[ \frac{1}{p^2} \sum_{j=1}^{N} \langle v_{a+b \bmod p}^{(w_j)} | v_c^{(w_j)} \rangle \quad \underset{\sim}{0} \quad \underset{\sim}{0} \right]$$

Projecting out the zero entries, the set of predictive vectors is given by:

$$V_p^{(w)} = \left\{ \langle v_i^{(w)}| \mid i = 0, 1, \ldots, p-1 \right\}$$

$\left[ \cos\left(\frac{2\pi w i}{p}\right) \quad \sin\left(\frac{2\pi w i}{p}\right) \right]$

vertices of a $p$-gon

## Summary:

| Process | Output | Predictive vectors |
|---|---|---|
| $RRMod_p$ | $P(w) = \langle \eta \mid T^{(w)} \mid 1 \rangle$ | $(p-1)$ — simplex |
| $sRRMod_p$ | $Z(w) = \langle \eta \mid H^{(w)} \mid 1 \rangle$ | $p$ — gon |

What do models represent?

# Results

## Recipe:

1) Train a one-layer one-head transformer to "grok" modular addition

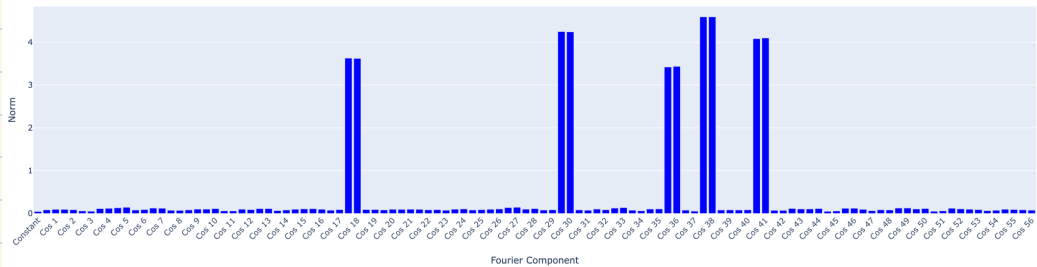2) Fit the logits of the transformer to:

$$\sum_{i=1}^{N} \cos\left(\frac{2\pi w_i (a+b-c)}{p}\right)$$

to determine $\underline{w} = (w_1, w_2, \ldots, w_N)$ of $s$ RR Mod $p$

3) Analyse model activations with:
   * linear regression
   * PCA

# Find frequencies: [Nanda et al.]



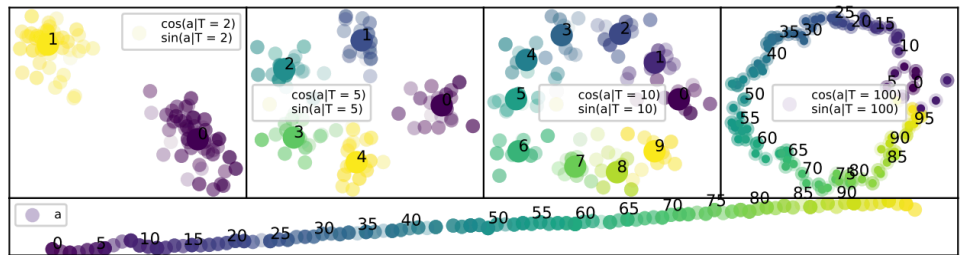Unembedding Fourier Component Norms - bos_single_head_cpu_20250718_144902 (Epoch 25000)

Simplex Linear Regression R² Evolution: [BOS, a, b] → (p-1)-Simplex Vertices
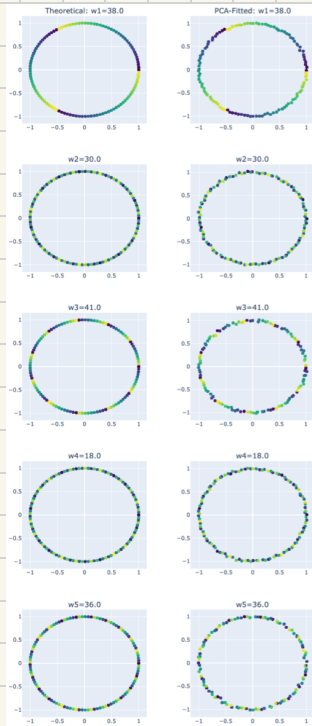p=113 | Target: One-hot vectors at position (a+b) mod p

Fourier Linear Regression R² Evolution: [BOS, a, b] → Fourier Components
p=113 | Frequencies: [38.0, 30.0, 41.0, 18.0, 36.0] | Target: cos/sin components of (a+b) mod p

Theoretical: w1=38.0 | PCA-Fitted: w1=38.0

w2=30.0 | w2=30.0

w3=41.0 | w3=41.0
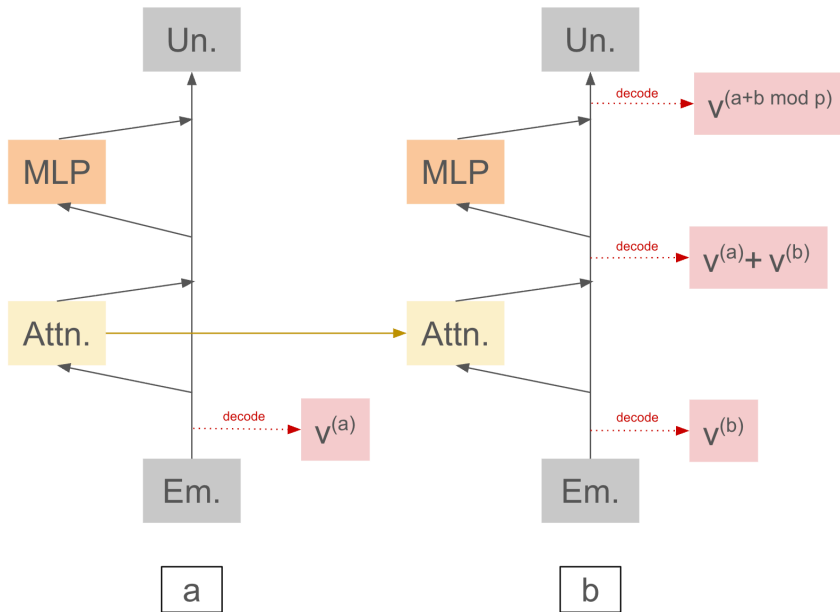
w4=18.0 | w4=18.0

w5=36.0 | w5=36.0

These results seem like toy versions of the [Kantamnemi et al.] results:
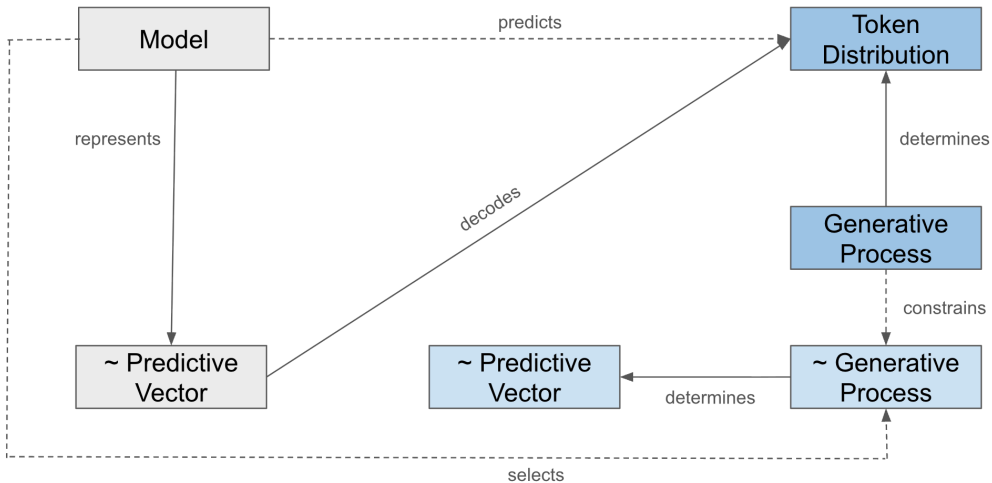
*Bonus !*



... actually represents more fourier components see [Yip et al.] for details

# Outlook

## Questions:

* What if we don't initialise in a synchronised state?

* If we directly train models to predict EHMM processes, are the predictive vectors decodable from activations?

* Is there a EHMM corresponding to familiar HMMs, e.g., is there an EHMM for Mess3?

# Thanks for Listening !