

# Singular Learning Theory 20 - In-context learning

sth20  
①  
25/11/22

The nature of in-context learning in Transformer models is likely to deeply inform any future theory which links the geometry of singularities in Transformers as learning machines to the emergent phenomena such as reasoning and abstraction which make these systems interesting.

In each layer of a Transformer model [PH] a list of entity representations  $e_1, \dots, e_\ell \in V$  are processed to form new representations  $e'_1, \dots, e'_\ell$ . This processing depends on the weights of the Transformer, in several ways

- (i) The layer normalisation
- (ii) The attention mechanism
- (iii) The feedforward layer in each block.

The final prediction also depends on embedding and unembedding weights, which we ignore here. We focus on (ii) the attention weights and how they are changed by backpropagation. The claim in [ICL1, ICL2] and other papers is that in-context learning is "similar" to this change in attention weights. These papers are not to be taken seriously, but let us investigate the question.

The cartoon Let  $T_w$  denote the Transformer model with attention weights  $w$ , as a family of functions  $\{T_w^{(n)}\}_{n \geq 1}$  where  $T_w^{(n)}$  takes a sequence of tokens  $\underline{x} = x_1, \dots, x_n$  as input and returns the prediction of the next token as a probability distribution. So the contribution of a sequence  $x_1, \dots, x_\ell, x_{\ell+1}$  to a gradient step is

$$\begin{aligned} \Delta w &= -\eta \nabla_w (-\log T_w^{(n)}(\underline{x})[x_{\ell+1}]) \\ &\propto \nabla_w T_w^{(n)}(\underline{x})[x_{\ell+1}] \end{aligned} \quad (1.1)$$

Suppose the tokens were examples of a task

$$\underbrace{x_1, \dots, x_l}_{\text{question}} \quad \underbrace{x_{l+1}}_{\text{answer}} \quad (2.1)$$

Then in-context learning refers to the possibility that on a second question  $y_1, \dots, y_l$

$$\overset{(2l+1)}{T_w(x, x_{l+1}, y)} \text{ is closer to } y_{l+1} \text{ than } \overset{(1)}{T_w(y)} \quad (2.2)$$

"learned from  $x$ "

There is of course another way to "improve"  $T_w$  given  $x, x_{l+1}$  which is to take the above gradient step. Which motivates an interesting question: is there some similarity between learning by gradient descent and in-context learning? More formally, is there some relationship between (for some  $\mathcal{Z}$ )

$$Q1. \quad T_w(x, x_{l+1}, -) \text{ and } T_{w+\Delta w}(-) ? \quad (2.3)$$

Providing a context is a form of currying where you can generate a new model, and the question is, in part, is this new model a Transformer for some perturbed weights? And if so in what direction?

$$Q2. \quad T_w(x, x_{l+1}, -) \approx T_{w'}(-) \quad \text{how does } w' \text{ relate to } w?$$

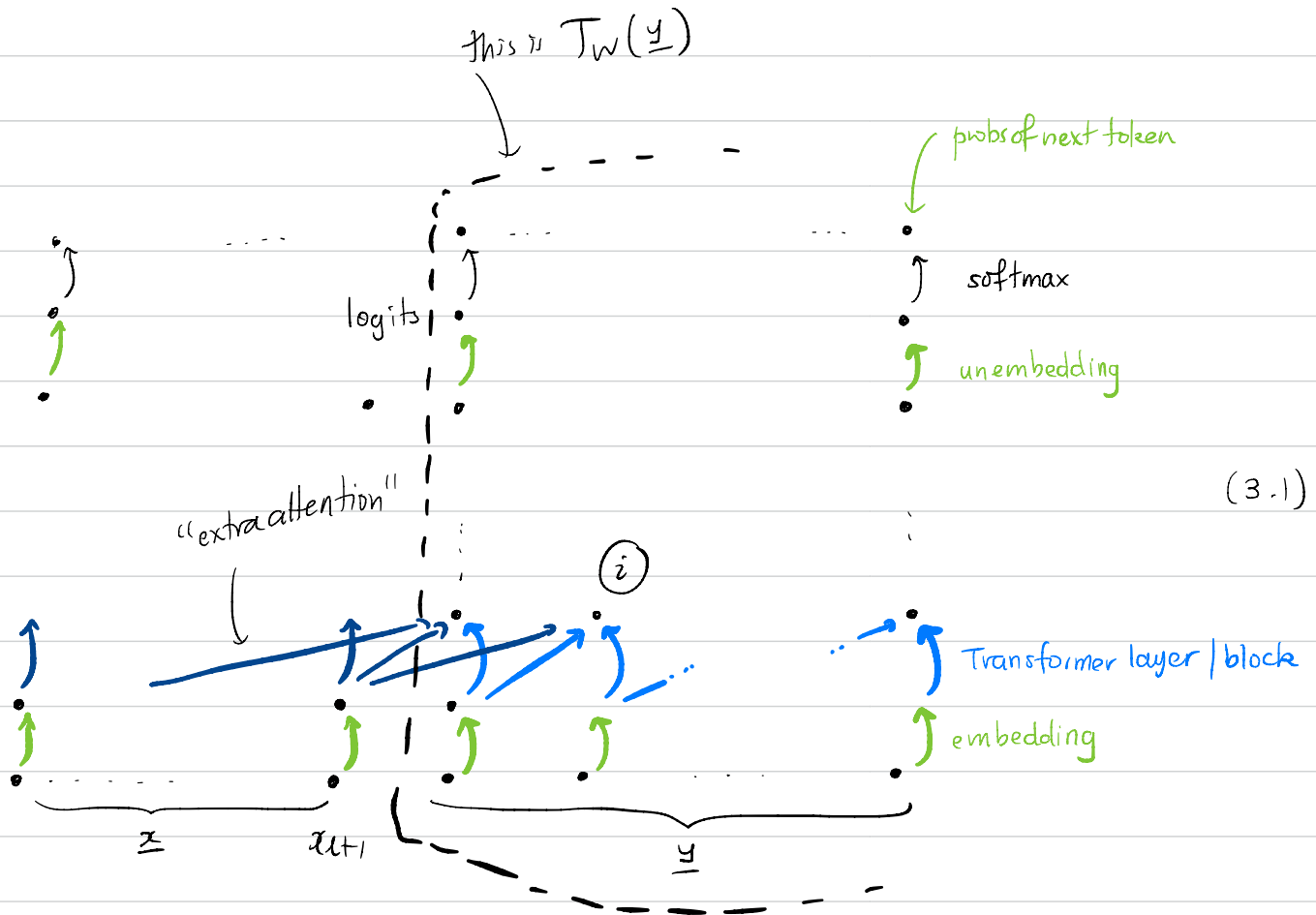
"contexts to weights hypothesis"

Quite unusually for statistical models, the attention weights actually specify an infinite family of models  $\{T_w^{(n)}\}_n$  and providing contexts allows us to move between them, potentially internalising many operations such as training which are currently on the outside.

Remark Contra [ICL1, ICL2] there seems little reason to believe in general that  $\|W - W'\|$  is small. For instance for the treacherous turn  $x =$  "In deployment".

We can probably safely say "No" to Q1, except in special cases. We focus on Q2.

Let us consider the algorithm computing  $T_w(\underline{x}, x_{t+1}, \underline{y})$



To simplify matters we are going to assume the entities  $e_1, \dots, e_{t+1}$  are frozen, and consider the effect on entity  $i$  of having  $\underline{x}, x_{t+1}$  in the context.

In the attention step the attention weight matrix  $W: V \rightarrow H \oplus H \oplus V$  generates queries, keys and value vectors (notation of [25]) (we ignore biases)

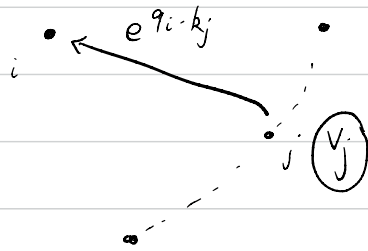
$$(\underbrace{q_i}_{in H}, \underbrace{k_i}_{in V}, v_i) := W e_i = (W^Q e_i, W^K e_i, W^V e_i)$$

and the update rule is (ignoring layer normalisation and feedforward)

$$e'_i = e_i + \sum_j \text{softmax}(q_i \cdot k_1, \dots, q_i \cdot k_e)_j v_j$$

$$= e_i + \sum_j \frac{e^{q_i \cdot k_j}}{\sum_u e^{q_i \cdot k_u}} v_j$$

$\Delta e_i$



If we define  $\text{softmax}(A)$  for a matrix  $A$  column-wise (again following [PH])

$$\text{softmax}(A)_{ij} = \frac{e^{A_{ij}}}{\sum_u e^{A_{uj}}}$$

Then  $\text{softmax}(K^T Q)_{ji} = \frac{e^{(K^T Q)_{ji}}}{\sum_u e^{(K^T Q)_{ui}}}$

$$= \frac{e^{q_i \cdot k_j}}{\sum_u e^{q_i \cdot k_u}}$$

$i$ th col  
↓

$$Q = W^Q E \quad E[:, i] = e_i$$

$$K = W^K E \quad Q[:, i] = q_i$$

$$e^{(K^T Q)_{ji}} = e^{q_i \cdot k_j}$$

Hence with  $E'[:, i] = e'_i$ , (1.2) becomes

$$E' = E + V \text{softmax}(K^T Q)$$

$H \leftarrow \text{ent} \leftarrow \text{ent}$

Returning to our situation where entities  $e_1, \dots, e_{t+1}$  are frozen

$$\Delta e_i = \sum_j \frac{e^{q_i \cdot k_j}}{\sum_u e^{q_i \cdot k_u}} v_j$$

$$= \frac{1}{\sum_u e^{q_i \cdot k_u}} \left\{ \sum_{j \text{ frozen}} e^{q_i \cdot k_j} v_j + \sum_{j \text{ not frozen}} e^{q_i \cdot k_j} v_j \right\}$$

Some special cases If  $q_i \cdot k_j \approx 0$  for frozen  $j$ , and  $q_i \cdot k_j$  sufficiently large for other  $j$ ,

$$\Delta e_i \approx \frac{1}{\sum_{u \text{ not frozen}} e^{q_i \cdot k_u}} \left\{ \sum_{j \text{ not frozen}} e^{q_i \cdot k_j} v_j \right\}$$

"Context is irrelevant"

$$= \text{update in } \mathcal{T}_w(\underline{y})$$

Suppose that for each frozen  $t$ , there is an unfrozen index  $a(t)$  such that  $v_j \approx v_{a(t)}$  and write

$$\Delta e_i \approx \frac{1}{\sum_u e^{q_i \cdot k_u}} \left\{ \sum_{\substack{j \text{ not} \\ \text{frozen}}} \left( e^{q_i \cdot k_j} + \sum_{\substack{t \text{ frozen} \\ a(t)=j}} e^{q_i \cdot k_t} \right) v_j \right\}$$

The Log-Sum-Exponential  $LSE(x_1, \dots, x_n) = \log(\sum_i e^{x_i})$  is an approximation to  $\max\{x_i\}$ : for  $x$ -values large relative to  $\log n$ , A common approx is  $(x^* = \max\{x_i\})$

$$LSE(x_1, \dots, x_n) \approx x^* + \log(\exp(x_1 - x^*) + \dots + \exp(x_n - x^*))$$

$$x^* \leq LSE \leq x^* + \log(n)$$

$$e^{q_i \cdot k_j} + \sum_{\substack{t \text{ frozen} \\ a(t)=j}} e^{q_i \cdot k_t} = e^z$$

$$z = \text{LSE}(\{q_i \cdot k_j\} \cup \{q_i \cdot k_t \mid t \text{ frozen } a(t)=j\})$$

$$\approx q_i \cdot k_{j^*} \quad \text{where this is maximised,}$$

Suppose there is a clear winner, for each unfrozen index  $j$ , then

$$\Delta e_i \approx \frac{1}{\sum_u e^{q_i \cdot k_u}} \left\{ \sum_{\substack{j \text{ not} \\ \text{frozen}}} e^{q_i \cdot k_{j^*}} v_j \right\}$$

which is attention but with modified key weights  $W^K$  for entity  $j$ .  
Perhaps sometimes the winner is a frozen entity, and in this way

$$T_w(\underline{x}, x_{t+1}, \underline{y}) \approx T_{w'}(\underline{y})$$

where  $W'$  includes the modified weights.

Remark • Treat properly using SVD and RA ideas.

## In-context learning and SLT

A learning machine in SLT is a tuple  $(\mathcal{W}, p, q, \mathcal{J})$  consisting of a parameterspace  $\mathcal{W}$ , a model  $p$ , truth  $q$  and prior  $\mathcal{J}$ . Let us formulate such a learning machine for a large language model like GPT. Note that "GPT" refers to more than just the Transformer model, it also refers to the data distribution and to some extent the training procedure.

Data distribution given some maximum context size  $C$ , sample sequences of consecutive tokens of length  $C$  from a fixed corpus of tokenstrings (e.g. text)

$$t_1, \dots, t_c \tag{7.1}$$

Then we have input-output pairs

$$\overset{(1)}{(t_1, t_2)}, (t_1 t_2, t_3), \dots, (t_1 \dots t_{c-1}, t_c) \tag{7.2}$$

The data distribution is the distribution of all examples like (7.2). So if  $S$  is the set of possible tokens, we have actually true distributions

$$\begin{aligned} q_1(x, y) & \text{ on } S \times S && (\text{prob. of } (t_1, t_2)) \\ q_2(x, y) & \text{ on } S^2 \times S && (\text{prob. of } (t_1 t_2, t_3)) \\ & \vdots && \\ q_{c-1}(x, y) & \text{ on } S^{c-1} \times S && \end{aligned}$$

Note that we can make predictions  $T_w(t_1), \dots, T_w(t_1 \dots t_{c-1})$  for all these examples with the same set of weights. The loss function for the Transformer is a sum over the cross-entropies for each of these predictions.

$$L(W) = \sum_{l=1}^{c-1} L^{(l)}(W) \quad (8.1)$$

where for  $1 \leq l \leq c-1$ ,

$$L^{(l)}(W) = \sum_{(\underline{t}, t') \in S^l \times S} q_l(\underline{t}) CE(q_l(t' | \underline{t}), T_w(\underline{t})) \quad (8.2)$$

Thus the true distribution is on  $X = (\prod_{l=1}^{c-1} S^l) \times S$  and given by  $q_1, \dots, q_{c-1}$  taken together. The model is, given  $\underline{t} \in S^l, t' \in S$ ,

$$\begin{aligned} p(\underline{t}, t' | \omega) &= p(t' | \underline{t}, \omega) q_l(\underline{t}) \\ &= T_w(\underline{t})[t'] q_l(\underline{t}) \end{aligned}$$

and the KL divergence is similarly a sum over  $l$

$$\begin{aligned} K(\omega) &= \int_x q(\underline{t}, t') \log \left( \frac{q(\underline{t}, t')}{p(\underline{t}, t' | \omega)} \right) d\underline{t} dt' \\ &= \int_x q(\underline{t}, t') \log \left( \frac{q(t' | \underline{t})}{p(t' | \underline{t}, \omega)} \right) d\underline{t} dt' \quad (8.3) \\ &= \int_x q(t' | \underline{t}) \log \left( \frac{q(t' | \underline{t})}{p(t' | \underline{t}, \omega)} \right) q(\underline{t}) d\underline{t} dt' \\ &= \sum_{l=1}^{c-1} \int_{S^l \times S} q(t' | \underline{t}) \log \left( \frac{q(t' | \underline{t})}{p(t' | \underline{t}, \omega)} \right) q_l(\underline{t}) d\underline{t} dt' \\ &=: \sum_{l=1}^{c-1} K^{(l)}(\omega) \end{aligned}$$

where  $K^{(l)}$  is the KL divergence for  $(W, p_l, q_l, \mathcal{P})$ .



By a familiar calculation

$$\begin{aligned} K^{(l)}(w) &= -\text{Entropy } q_{\mathcal{L}} + \int_{S^{\mathcal{L}}} \text{CE}(q_{\mathcal{L}}(t'/\underline{t}) \mid T_w(\underline{t})) q_{\mathcal{L}}(\underline{t}) d\underline{t} \\ &= -\text{Entropy } q_{\mathcal{L}} + L^{(l)}(w) \end{aligned} \quad (9.1)$$

so minimising  $K^{(l)}(w)$  is the same as minimising  $L^{(l)}(w)$ .

Remark Since the same parameter  $w \in \mathcal{W}$  is used to make next-token predictions for any length this is not a sum of KL divergences for distinct models.

### The Bayesian posterior of a language model

The nature of the data distribution (that we see, for a sequence of tokens  $x_1, \dots, x_n$  every initial segment  $x_1, \dots, x_n$ ) and of the model (we use the same parameters to predict for any sequence length) have some interesting implications for the Bayesian posterior when we incorporate the contexts-to-weights hypothesis.

Recall that given a set of samples  $D_n = \{(\underline{t}^{(i)}, t'^{(i)})\}_{i=1}^n$  the posterior is

$$\begin{aligned} p(w|D_n) &= \frac{p(D_n|w)p(w)}{p(D_n)} \\ &= \frac{1}{Z_n} \mathcal{J}(w) \exp(-nL_n(w)) \end{aligned} \quad (9.2)$$

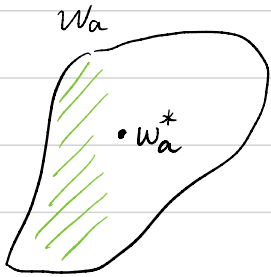
where  $L_n(w)$  is the empirical loss

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \text{CE}(q_{\mathcal{L}}(t'^{(i)} \mid \underline{t}^{(i)}), T_w(\underline{t}^{(i)})) \\ = -\frac{1}{n} \sum_{i=1}^n \log T_w(\underline{t}^{(i)})[t'^{(i)}] \end{aligned} \quad (9.3)$$

The model evidence  $Z_n = \int dw \mathcal{P}(w) \exp(-n L_n(w))$  is the basis of model selection in a Bayesian framework. We view model selection as internalized in LLMs to phase selection where two regions  $\mathcal{W}_1, \mathcal{W}_2 \subseteq \mathcal{W}$  are preferred in accordance with which has higher "evidence"

$$Z_n(\mathcal{W}_a) := \int_{\mathcal{W}_a} dw \mathcal{P}(w) \exp(-n L_n(w)) \quad a \in \{1, 2\}$$

Or what is the same, which has lower free energy  $F_a := -\log Z_n(\mathcal{W}_a)$  (the effective Boltzmann weight). Under some conditions we have the



Free Energy Formula

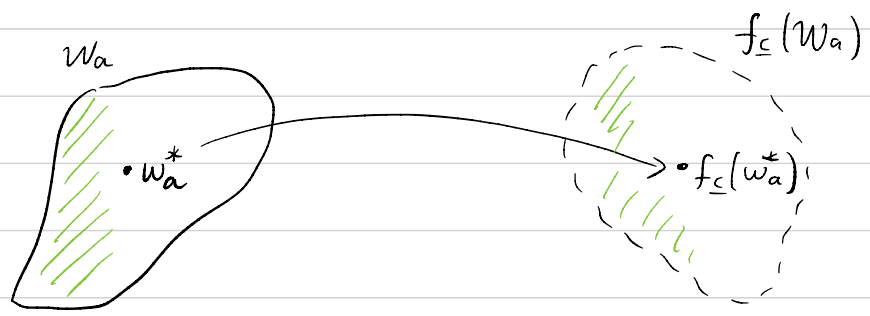
$$Z_n(\mathcal{W}_a) \approx n L_n(w_a^*) + \lambda_a \log n$$

(leading terms in asymptotic exp. see Green book, WBLE paper, DLT3)

Each point  $(\underline{t}^{(i)}, t'^{(i)})$  contributes to free energy of every phase.

$$\begin{aligned} Z_n(\mathcal{W}_a) &= \int_{\mathcal{W}_a} dw \mathcal{P}(w) \exp\left(\sum_{i=1}^n \log T_w(\underline{t}^{(i)})[t'^{(i)}]\right) \\ &= \int_{\mathcal{W}_a} dw \mathcal{P}(w) \prod_{i=1}^n T_w(\underline{t}^{(i)})[t'^{(i)}] \end{aligned}$$

Suppose The contexts-to-weights hypothesis holds for some sample  $D_n$  as above, with  $\underline{t}^{(i)} = \underline{c} \underline{x}^{(i)}$  for all  $i$  (in practice we find in  $D_n$  some subset like this) in the sense that  $T_w(\underline{c} \underline{x}^{(i)})[t'^{(i)}] \approx T_{f_{\underline{c}}(w)}(\underline{x}^{(i)})[t'^{(i)}]$  and moreover assume this works for all  $w \in \mathcal{W}_a$  in some way continuous in  $w$ . That is



Note  $f_c$  is likely to be highly degenerate

$$\begin{aligned}
 Z_n(W_a, D_n) &= \int_{W_a} dw \mathcal{P}(w) \prod_{i=1}^n T_w(\underline{t}^{(i)}) [t'^{(i)}] \\
 \text{"few shot"} &= \int_{W_a} dw \mathcal{P}(w) \prod_{i=1}^n T_w(\leq \underline{x}^{(i)}) [t'^{(i)}] \\
 &\approx \int_{W_a} dw \mathcal{P}(w) \prod_{i=1}^n T_{f_c(w)}(\underline{x}^{(i)}) [t'^{(i)}] \\
 &= \int_{f_c(W_a)} dw \bar{\mathcal{P}}(w) \prod_{i=1}^n T_w(\underline{x}^{(i)}) [t'^{(i)}] \\
 &\quad \uparrow \text{absorbs } \mathcal{J} \\
 &< Z_n(f_c(W_a), (\underline{x}^{(i)}, t'^{(i)})_{i=1}^n) \\
 &\quad \uparrow \text{"zero shot"}
 \end{aligned}$$

Since in the asymptotic expansion  $f_c$  contributes additional degeneracy on the LHS. Hence

$$n \ln(w_a^*) + \lambda \log n < n \ln(f_c(w_a^*)) + \lambda_a^c \log n$$

—  $\lambda_a^c = \lambda_a + \text{"degeneracy of } f_c \text{"}$

and so (the rough argument goes)  $\lambda_a < \lambda_a^c$ . The singularity at  $w_a^*$  must be more complex than the one at  $f_c(w_a^*)$  (roughly because  $f_c(w_a^*)$  is "specialised" and this corresponds to  $f_c$  being a projection, i.e. being degenerate. The more specialised the model becomes in context  $c$ , the larger the gap).

Phase as UTM, contexts as codes

Suppose  $W_a$  knows to direct predictions to other phases (subroutines) for multiple contexts  $\underline{\epsilon}_1, \underline{\epsilon}_2$ . So we suppose a dataset  $D_n$  consists of

$$D_n = \left\{ (\underline{\epsilon}_1, \underline{x}^{(i)}, y^{(i)}) \right\}_{i=1}^{n_1} \cup \left\{ (\underline{\epsilon}_2, \underline{x}^{(j)}, y^{(j)}) \right\}_{j=1}^{n_2} \quad n_1 + n_2 = n$$

Then

$$\begin{aligned} Z_n(W_a, D_n) &= \int_{W_a} d\omega \mathcal{P}(\omega) \prod_{i=1}^{n_1} T_\omega(\underline{\epsilon}_1, \underline{x}^{(i)}) [y^{(i)}] \\ &\quad \cdot \prod_{j=1}^{n_2} T_\omega(\underline{\epsilon}_2, \underline{x}^{(j)}) [y^{(j)}] \\ &\approx \int_{W_a} d\omega \mathcal{P}(\omega) \prod_{i=1}^{n_1} T_{f_{\underline{\epsilon}_1}(\omega)}(\underline{x}^{(i)}) [y^{(i)}] \\ &\quad \cdot \prod_{j=1}^{n_2} T_{f_{\underline{\epsilon}_2}(\omega)}(\underline{x}^{(j)}) [y^{(j)}] \end{aligned}$$

Suppose  $f_{\underline{\epsilon}_1}$  is constant in all the directions  $f_{\underline{\epsilon}_2}$  has nonzero partial derivatives at  $w_a^*$  and vice-versa, so there are local coordinates  $u, v$  in which  $f_{\underline{\epsilon}_1}(u, v) = f_{\underline{\epsilon}_1}(u)$ ,  $f_{\underline{\epsilon}_2}(u, v) = f_{\underline{\epsilon}_2}(v)$ . Then (if this is the only degeneracy)

$$\begin{aligned} &\approx \int_{W_a} du dv \mathcal{P}(u, v) \prod_i T_{f_{\underline{\epsilon}_1}(u)}(\dots) \prod_j T_{f_{\underline{\epsilon}_2}(v)}(\dots) \\ &= \int du \mathcal{P}(u) \prod_i T_{f_{\underline{\epsilon}_1}(u)}(\dots) \int dv \mathcal{P}(v) \prod_j T_{f_{\underline{\epsilon}_2}(v)}(\dots) \end{aligned}$$

[w, Remark 2.2]

$$-\log Z_n \approx n \ln(w_a^*) + (\lambda_a^{\underline{\epsilon}_1} + \lambda_a^{\underline{\epsilon}_2}) \log n$$

We see here a hint of a picture where the true distribution can be approximately divided into subtasks, each of which is much simpler in isolation than the true distribution, and under the context-to-weights hypothesis some contributions to the free energy of the "mother singularity" or phase  $W_a$  come from the free energy of submodels ( $f_{\subseteq}(W_a)$  lower-dim) on these tasks.

This suggests an approach to understanding the phase structure of the full model in terms of

- (A) Phase transitions where a given context is recognised  
(i.e. the context-to-weights hypothesis begins to be applied)
- (B) Phase transitions in the submodels
- (C) Interactions between the transitions in (B). In the generic case they are simply superimposed.

## References

- [PH] M. Phuong, M. Hutter "Formal Algorithms for Transformers" 2022.
- [GPT3] T. Brown et al "Language models are few shot learners" 2020.
- [ICL1] J. van Osvald et al "Transformers learn in-context by gradient descent" 2022
- [ICL2] D. Dai et al "Why can GPT learn in context?" 2022.
- [25] J. Clift et al "Logic and the 2-simplicial Transformer" ICLR 2020