Is Attention Better Than Matrix Decomposition?

Zhengyang Geng^{1,2}, Meng-Hao Guo^{3,*},

Hongxu Chen⁴, Xia Li², Ke Wei⁴, Zhouchen Lin^{2,5,\$}

¹Zhejiang Lab; ²Key Lab. of Machine Perception (MoE), School of EECS, Peking University; ³Tsinghua University^{; 4}School of Data Science, Fudan University^{; 5}Pazhou Lab; *Equal First Authorship; ^{\$}Corresponding Author

ICLR 2021







1. When we talk about global information/long-range correlation in attention-related methods, what do we actually mean?

Model the global information in a first-principle way.

2. Is hand-crafted attention irreplaceable when modeling the global context?

Provide a strong white-box baseline Hamburger for attention mechanism.







Matrix Decomposition





decomposition

 $\mathrm{rank}(oldsymbol{\bar{X}}) \leq \min(\mathrm{rank}(oldsymbol{D}),\mathrm{rank}(oldsymbol{C})) \leq r \ll \min(d,n).$

1. Inverse of Generation

- 2. Low-Rankness
- 3. Optimization

Hamburger





$$\min_{oldsymbol{D},oldsymbol{C}} ~~ \mathcal{L}(oldsymbol{X},oldsymbol{D}oldsymbol{C}) + \mathcal{R}_1(oldsymbol{D}) + \mathcal{R}_2(oldsymbol{C})$$

- Learn the global information <=> Solve the the optimization of MD
- Denote the optimization algorithm to solve the problem as $\mathcal{M}.\mathcal{M}$ is the core architection of Hamburger.

Ham



VQ

NMF

 $\min_{oldsymbol{D},oldsymbol{C}} \|oldsymbol{X} - oldsymbol{D}oldsymbol{C}\|_F \hspace{1em} ext{s. t. } \mathbf{c}_i \in \{\mathbf{e}_1, \mathbf{e}_2, \cdots, \mathbf{e}_r\} \hspace{1em} \max_{oldsymbol{D}}$

$$\min_{oldsymbol{D},oldsymbol{C}} \|oldsymbol{X} - oldsymbol{D}oldsymbol{C}\|_F \quad ext{s. t. } oldsymbol{D}_{ij} \geq 0, oldsymbol{C}_{jk} \geq 0$$

Algorithm 1 Ham: Soft VQ	Algorithm 2 Ham: NMF with MU
Input X . Initialize D, C .	Input X . Initialize non-negative D, C
for k from 1 to K do	for k from 1 to K_do
$oldsymbol{C} \leftarrow softmax(rac{1}{T}cosine(oldsymbol{D},oldsymbol{X}))$	$oldsymbol{C}_{ij} \leftarrow oldsymbol{C}_{ij} rac{(oldsymbol{D}^{ op}oldsymbol{X})_{ij}}{(oldsymbol{D}^{ op}oldsymbol{D}oldsymbol{C})_{ij}}$
$oldsymbol{D} \leftarrow oldsymbol{X}oldsymbol{C}^ op diag(oldsymbol{C}oldsymbol{1}_n)^{-1}$	$oldsymbol{D}_{ij} \leftarrow oldsymbol{D}_{ij} rac{(oldsymbol{X}oldsymbol{C}^ op)_{ij}}{(oldsymbol{D}oldsymbol{C}oldsymbol{C}^ op)_{ij}}$
end for	end for
Output $\bar{X} = DC$.	Output $\bar{X} = DC$.

One-step Grad

• We build an abstract model to analyse the grad back through \mathcal{M} .

 $\mathbf{h}^{t+1} = \mathcal{F}(\mathbf{h}^t, \mathbf{x})$

- We find the thorny **scale** and **spectrum** of the Jacobian matrix in the BPTT algorithm.
- Back-propagation through optimization suffers from vanishing gradient w.r.t. the initialization and exploding gradient w.r.t. the input!
- We use the one-step grad $\frac{\partial y}{\partial x}$, i.e., the grad from the last optimization step.





Table 1: One-step Gradient & BPTT

Method	One-step	BPTT
VQ	77.7(77.4)	76.6(76.3)
CD	78.1(77.5)	75.0(74.6)
NMF	78.3(77.8)	77.4(77.0)

A Close Look at Hamburger





Accumulative Ratio

Visualization of Feature Maps

Computation & Memory



Method	Params	MACs GPU Load		Load	GPU Time	
u	i ui uiiis	ivine 5	Train	Infer	Train	Infer
SA	1.00M	292G	5253MB	2148MB	242.0ms	82.2ms
DA	4.82M	79.5G	2395MB	2203MB	72.6ms	64.4ms
A^2	1.01M	25.7G	326MB	165MB	22.9ms	8.0ms
APC	2.03M	17.6G	458MB	264MB	26.5ms	11.6ms
DM	3.00M	35.1G	557MB	268MB	65.7ms	23.3ms
ACF	0.75M	79.5G	1380MB	627MB	71.0ms	22.6ms
Ham (CD)	0.50M	16.2G	162MB	102MB	20.0ms	13.0ms
Ham (NMF)	0.50M	17.6G	202MB	98MB	15.6ms	7.7ms

Table 3: Comparisons between Hamburger and context modules.

SOTA on Semantic Segmentation

Large Scale Image Generation

Table 4: Comparisons with state-of-the-art on the PASCAL VOC test set w/o COCO pretraining.

Method	mIoU(%)
PSPNet (Zhao et al., 2017)	82.6
DFN* (Yu et al., 2018)	82.7
EncNet (Zhang et al., 2018)	82.9
DANet* (Fu et al., 2019)	82.6
DMNet* (He et al., 2019a)	84.4
APCNet* (He et al., 2019b)	84.2
CFNet* (Zhang et al., 2019b)	84.2
SpyGR* (Li et al., 2020)	84.2
SANet* (Zhong et al., 2020)	83.2
OCR* (Yuan et al., 2020)	84.3
HamNet	85.9

Table 5: Results on the PASCAL-Context Val set.

Method mIoU(%) 47.8 PSPNet (Zhao et al., 2017) SGR* (Liang et al., 2018) 50.8 51.7 EncNet (Zhang et al., 2018) DANet* (Fu et al., 2019) 52.6 53.1 EMANet* (Li et al., 2019) DMNet* (He et al., 2019a) 54.4 APCNet* (He et al., 2019b) 54.7 CFNet* (Zhang et al., 2019b) 54.0 SpyGR* (Li et al., 2020) 52.8 SANet* (Zhong et al., 2020) 53.0 OCR* (Yuan et al., 2020) 54.8 HamNet 55.2

Table 6: Results on ImageNet 128×128 . * are from Tab. 1 and Tab. 2 of Zhang et al. (2019a).

Method	FID↓
SNGAN-projection*	27.62
SAGAN*	18.28
HamGAN-baby	16.05
YLG	15.94
HamGAN-strong	14.77



Performance





Global Information -> Low-rank Formulation

Optimization as Architecture

The Curse -> Gradient back through MD -> One-step Grad

Enjoy Hamburger!



Thank you!

