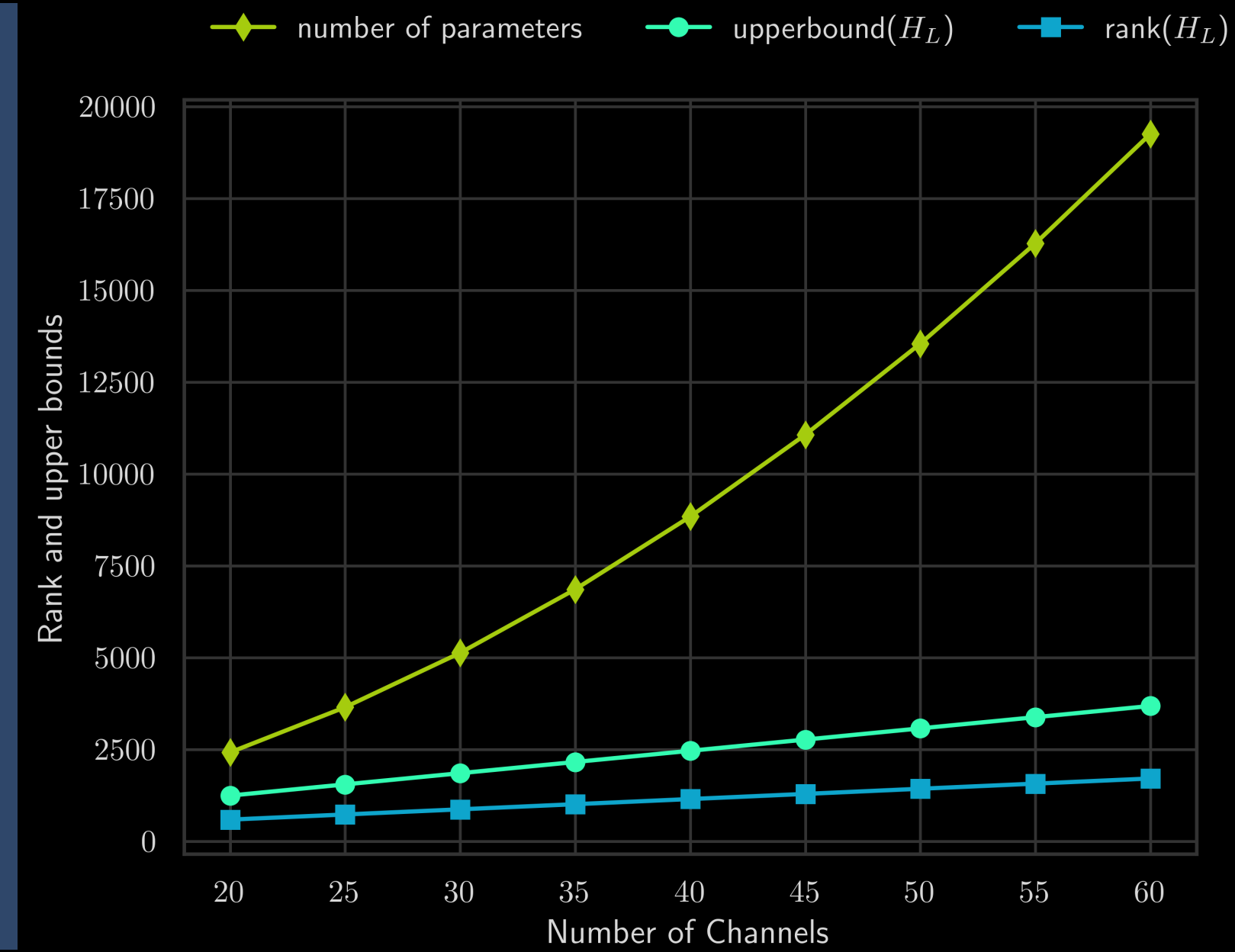


The Hessian Perspective into the Nature of Convolutional Neural Networks

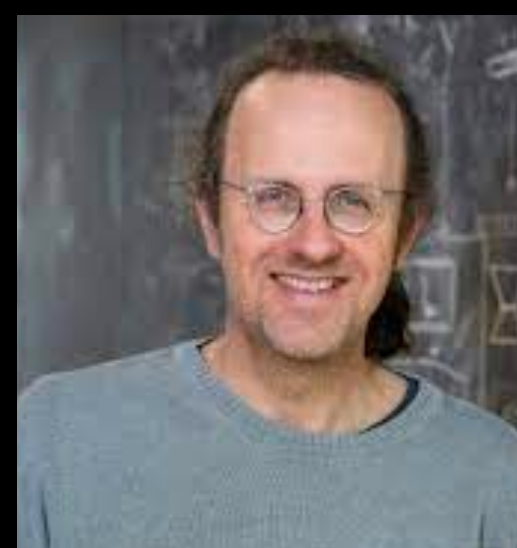
ICML 2023



Sidak Pal Singh



Thomas Hofmann



Bernhard Schölkopf

ETH zürich



Aim of the Work

To provide a perspective into the nature of CNNs, i.e., how their architectural characteristics of CNNs manifest themselves in terms of the properties of its loss landscape as given by the Hessian rank

Hessian Rank, and a question

Captures pairwise interactions of parameters (θ_i, θ_j)

via the second-derivatives of the loss \mathcal{L}

$$\mathbf{H}_{ij} = \frac{\partial^2 \mathcal{L}}{\partial \theta_i \partial \theta_j}$$

Hessian Range $\text{range}(\mathbf{H}) = \{y = \mathbf{H}\theta : \theta \in \mathbb{R}^p\}$

Hessian Rank $\text{rank}(\mathbf{H}) = \dim(\text{range}(\mathbf{H}))$

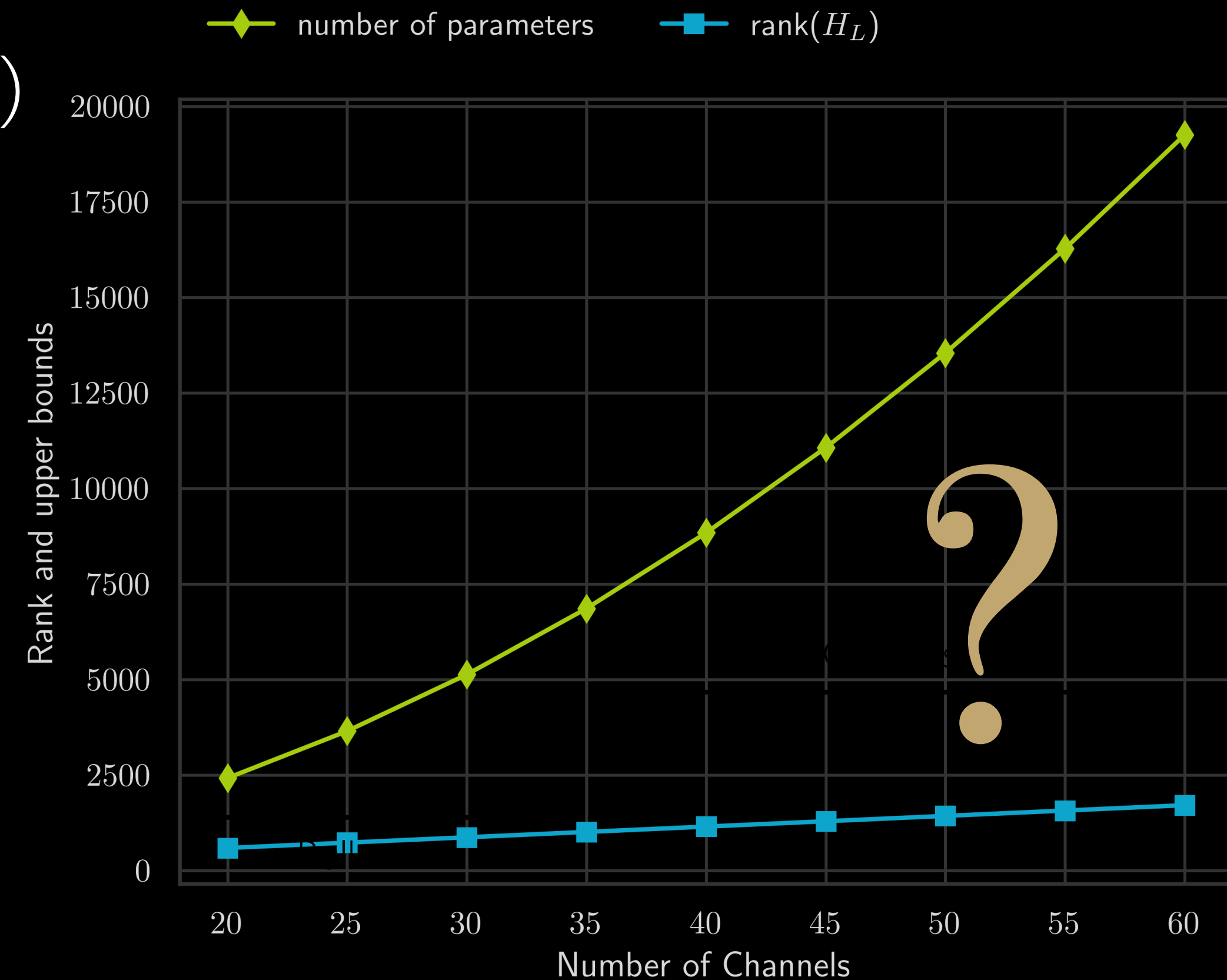
Hessian Rank, and a question

Captures pairwise interactions of parameters (θ_i, θ_j)
via the second-derivatives of the loss \mathcal{L}

$$\mathbf{H}_{ij} = \frac{\partial^2 \mathcal{L}}{\partial \theta_i \partial \theta_j}$$

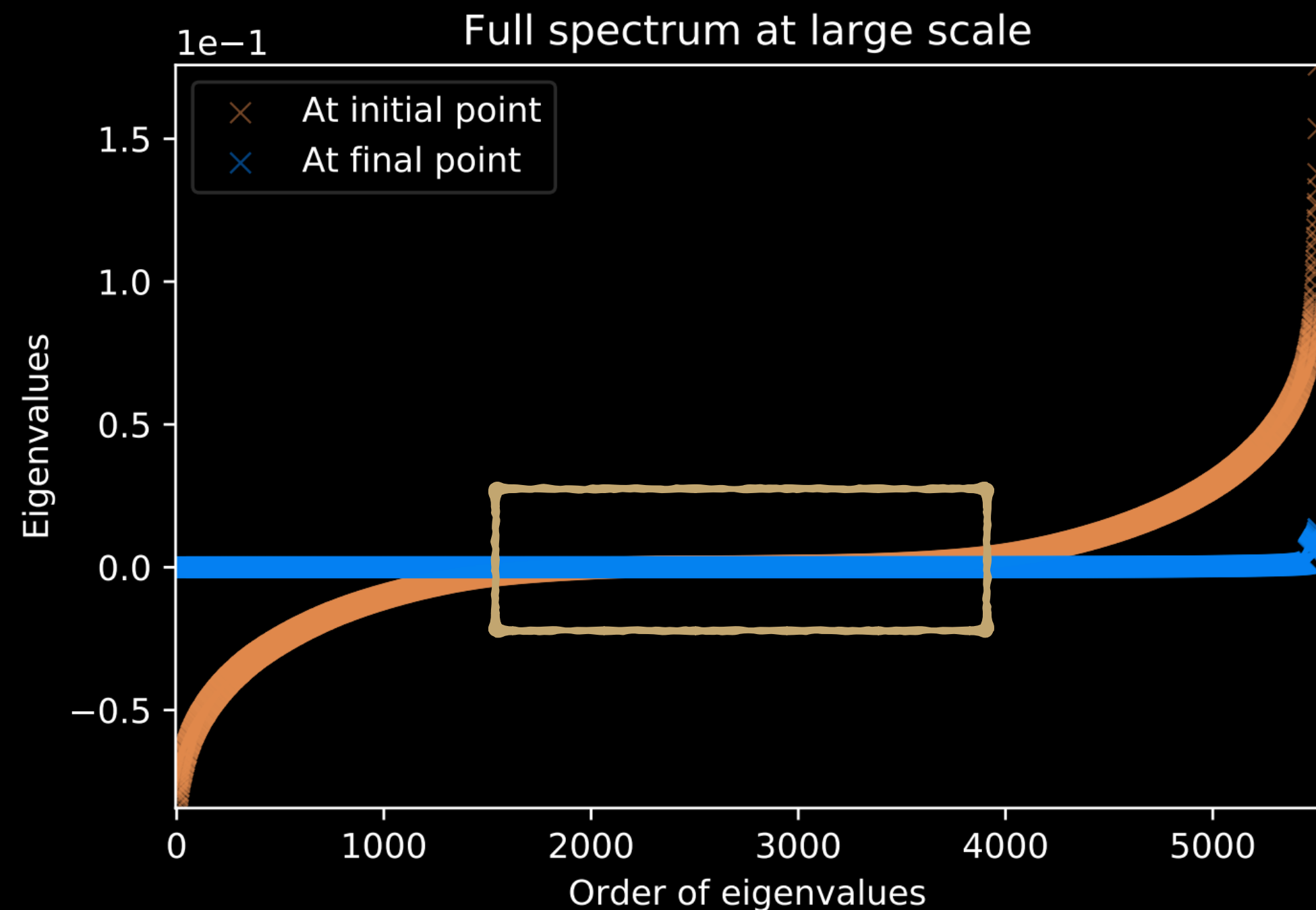
Hessian Range $\text{range}(\mathbf{H}) = \{y = \mathbf{H}\theta : \theta \in \mathbb{R}^p\}$

Hessian Rank $\text{rank}(\mathbf{H}) = \dim(\text{range}(\mathbf{H}))$



Related Work

Sagun et. al., 2017



Significant extent of degeneracy in the Hessian

Singh et. al., 2021

$$\text{rank}(\mathbf{H}_o) = q(d + K - q)$$

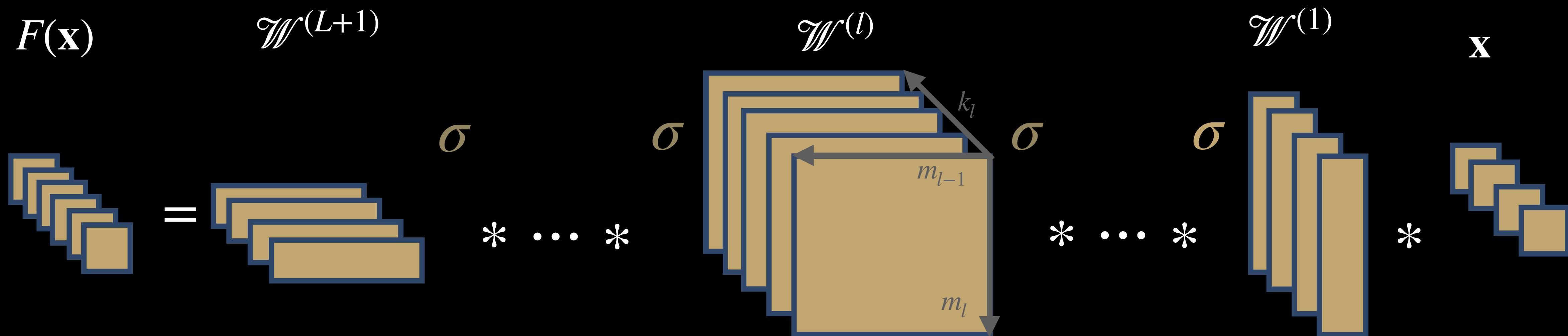
$$\text{rank}(\mathbf{H}_f) \leq 2q \sum_{i=1}^L m_i + 2qs - (L + 1)q^2$$

$$\text{rank}(\mathbf{H}_{\mathcal{L}}) = 2q \sum_{i=1}^L m_i - (L + 1)q^2 + q(r + K)$$

Theoretical characterisation for *Linear* Fully-Connected networks (FCNs)

Setup and Formalism

Given input $\mathbf{x} \in \mathbb{R}^{d_0}$, a deep CNN F with L hidden layers of weight tensors $\mathcal{W}^l \in \mathbb{R}^{m_l \times m_{l-1} \times k_l}$, nonlinearity $\sigma(\cdot)$, stride 1, zero padding, with input channels $m_0 = 1$, output channels = K



Shorthands: $\mathbf{T}^{(k:l)} := \mathbf{T}^{(k)} \dots \mathbf{T}^{(l)} \quad \forall k > l$ and $\mathbf{T}^{(k:l)} := \mathbf{T}^{(k)\top} \dots \mathbf{T}^{(l)\top} \quad \forall k < l$

Also, denote the set of all parameters by $\theta := \{\mathcal{W}^1, \dots, \mathcal{W}^{L+1}\}$

Toeplitz Framework: full-fledged CNNs

$$F_{\theta}(\mathbf{x}) = \mathcal{W}^{(L+1)} * \sigma \left(\mathcal{W}^{(L)} * \sigma \left(\dots * \sigma \left(\mathcal{W}^{(1)} * \mathbf{x} \right) \right) \right)$$

$$\mathbf{T}^{(l)} := \begin{pmatrix} \mathbf{T}^{\mathcal{W}^{(l)}_{(1,1)}} \cdot & \dots & \mathbf{T}^{\mathcal{W}^{(l)}_{(1,m_{l-1})}} \cdot \\ \vdots & & \vdots \\ \mathbf{T}^{\mathcal{W}^{(l)}_{(m_l,1)}} \cdot & \dots & \mathbf{T}^{\mathcal{W}^{(l)}_{(m_l,m_{l-1})}} \cdot \end{pmatrix}$$



$$F_{\theta}(\mathbf{x}) = \mathbf{T}^{(L+1)} \Lambda_{\mathbf{x}}^{(L)} \mathbf{T}^{(L)} \dots \Lambda_{\mathbf{x}}^{(1)} \mathbf{T}^{(1)} \mathbf{x}$$

$\Lambda_{\mathbf{x}}^{(l)}$ contains the activations at layer l

Each of the $\mathbf{T}^{\mathcal{W}^{(l)}_{(i,j)}}$ is a Toeplitz matrix as discussed before

Key theoretical results

Theorem 1: Rank of outer-product Hessian $\mathbf{H}_o = \mathbf{E}_p \left[\nabla_{\theta} F(\mathbf{x}) \nabla_{\theta} F(\mathbf{x})^{\top} \right]$

The rank of \mathbf{H}_o for a deep linear CNN, with kernel sizes k_l and number of filters m_l is:

$$\begin{aligned} \text{rank}(\mathbf{H}_o) &\leq \min \left(p, d_0 \text{rank}(\mathbf{T}^{(2:L+1)}) + K \text{rank}(\mathbf{T}^{(L:1)}) - \text{rank}(\mathbf{T}^{(2:L+1)}) \text{rank}(\mathbf{T}^{(L:1)}) \right) \\ &= \min(p, q_o(d_0 + K - q_o)). \end{aligned}$$

where $q_o := \min(d_0, m_1 d_1, \dots, m_L d_L, K)$ denotes the flattened bottleneck dimension.

- If there is no bottleneck within, $\text{rank}(\mathbf{H}_o) = Kd_0$

Key theoretical results

Theorem 2: Rank of functional Hessian $\mathbf{H}_f = \mathbf{E}_p \left[\sum_{c=1} [\partial \ell_{\mathbf{x}, \mathbf{y}}]_c \nabla_{\boldsymbol{\theta}}^2 F_c(\mathbf{x}) \right]$

The rank of the l -th column block of \mathbf{H}_f for deep linear CNN, with kernel sizes k_l and number of filters m_l is: $\text{rank}(\mathbf{H}_f^{\bullet l}) \leq \min(q_f m_{l-1} d_{l-1} + q_f m_l d_l - q_f^2, m_l m_{l-1} k_l),$

for $l \in [2, \dots, L]$ and where $q_f := \min(q_o, s)$ and $s := \text{rank}(\boldsymbol{\Omega}) = \text{rank}(\mathbf{E}[\boldsymbol{\delta}_{\mathbf{x}, \mathbf{y}} \mathbf{x}^\top])$.

Hence we have that, $\text{rank}(\mathbf{H}_f) \leq \sum_{l=1}^{L+1} \text{rank}(\mathbf{H}_f^{\bullet l})$ (eq. 1)

- **Block-column independence:** Like in the case of FCNs, simply adding the ranks of the block-columns of the respective layers, gives the rank of the entire \mathbf{H}_f without introducing any looseness in the bound (so, the inequality in eq. 1 suffices)

Key theoretical results

Rank of the loss Hessian can be bounded as

$$\text{rank}(\mathbf{H}_{\mathcal{L}}) = \text{rank}(\mathbf{H}_o + \mathbf{H}_f) \leq \text{rank}(\mathbf{H}_o) + \text{rank}(\mathbf{H}_f)$$

Rank of the loss Hessian grows as

$$\mathcal{O}(m \cdot L \cdot d_0)$$

Number of parameters grow as

$$\mathcal{O}(m^2 \cdot L \cdot d_0)$$

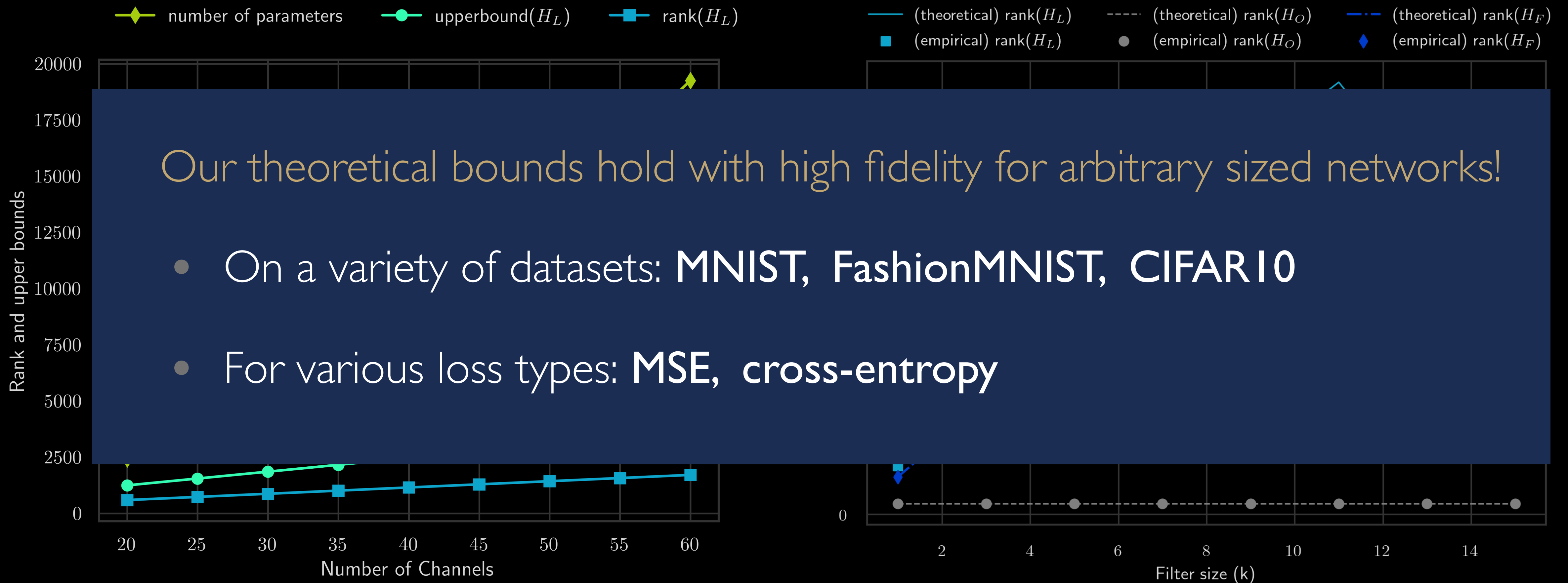
For typical networks, $m \gg L$ and $m \gg d_0$

Hence, rank will show a *square root behaviour* relative to the number of parameters

Thus generalizing the finding of Singh et. al. (2021) to the case of CNNs

Empirical Results

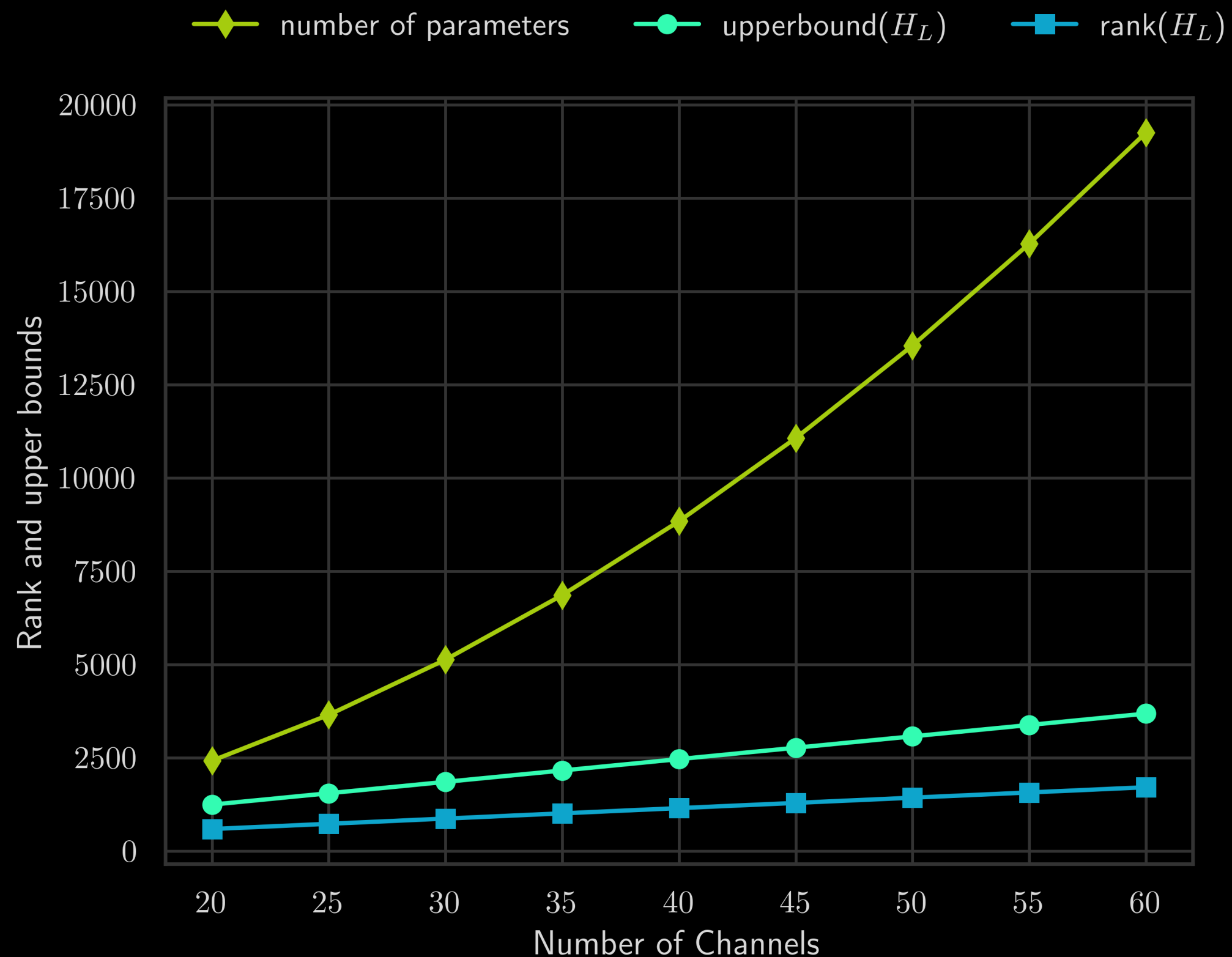
$$\min(\text{kernel-size}_i, \text{spatial-dim}_i)$$



Rank bounds for *increasing # of filters m*

Rank bounds for *increasing filter size k*

Summary



- Employ an equivalent representation of CNNs as composition of **Toeplitz** maps
- Natural change: $m_i \rightarrow m_i d_i$ where
$$d_i = \frac{d_{i-1} - \text{kernel-size}_i + 2 \text{padding}_i}{\text{stride}_i} + 1$$
- The **square-root** trend of rank persists

This sheds a novel perspective on the nature of CNNs and highlights the degree of redundancy inherent in over-parameterized networks.