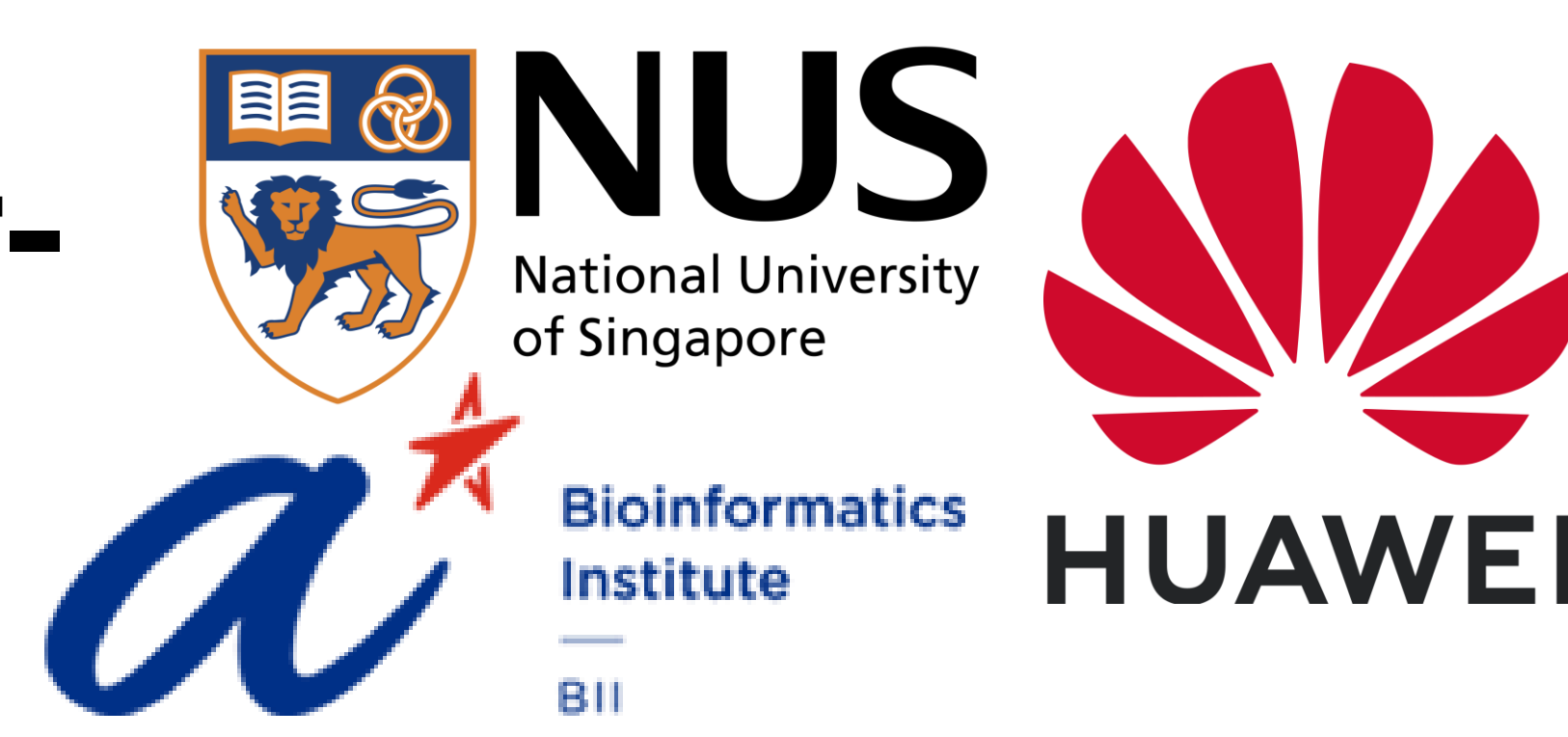


# In-Context Learning behaves as a greedy layer-wise gradient descent algorithm

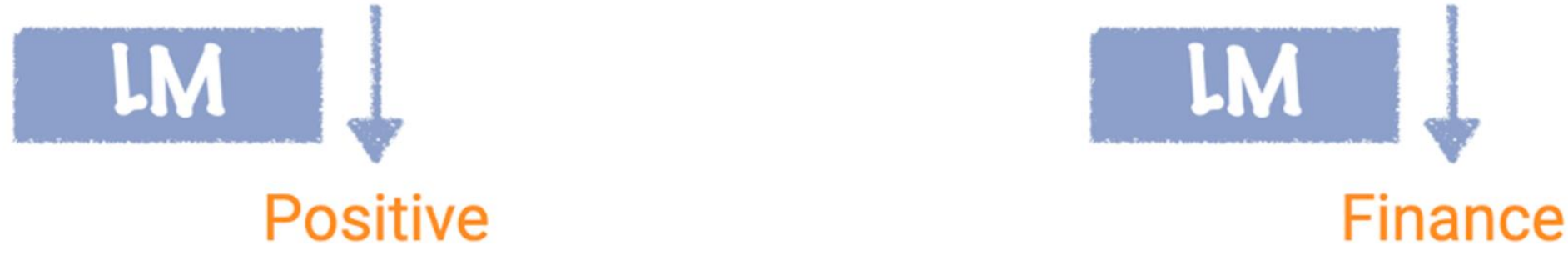
Brian K Chen, Tianyang Hu, Hui Jin, Hwee Kuan Lee, Kenji Kawaguchi



## Motivation

- ICL is an emergent property found in LLMs
- Appending prompts to existing ones “teaches” the LLM new information without training the model

Circulation revenue has increased by 5% in Finland. // Positive	Circulation revenue has increased by 5% in Finland. // Finance
Panostaja did not disclose the purchase price. // Neutral	They defeated ... in the NFC Championship Game. // Sports
Paying off the national debt will be extremely painful. // Negative	Apple ... development of in-house chips. // Tech
The company anticipated its operating profit to improve. // _____	The company anticipated its operating profit to improve. // _____



- Theoretical understanding of ICL remains limited
- Existing work that considers ICL as a single step of gradient descent focuses on limited cases with specific weights
- Want to study the mechanics of ICL by looking at the linearized attention module with generic weights

## Main Result

### Theorem 1:

For an initial self-attention mechanism with query token  $q$ , matrices  $W_V, W_K$ , prompt  $X = [p_N, \dots, p_1]$  and operator  $\mathcal{F}_0([X], q) = \text{LinAttn}(W_V X, \phi(W_K X), q)$ , the following systems are equivalent:

$$S_0 = \mathcal{F}_0([X'; X], q)$$

and  $S_1 = \mathcal{F}_1([X], q)$

Where  $\mathcal{F}_1([X], \cdot)$  is the linear function  $\mathcal{F}_0([X], \cdot)$  after one step of gradient descent with learning rate  $\eta$  and training set  $\{x_i, y_i\}_{i=1}^M$ . For every  $i \in \{1, \dots, M\}$ :

$$x_i = \phi(W_K p'_i)$$

$$y_i = \frac{M}{\eta} W_V p'_i + \mathcal{F}_0([X], \phi(W_K p'_i))$$

## Preliminaries

Linearization of the attention mechanism:

$$\text{LinAttn}(V, \phi(K), \phi(Q)) = V \phi(K)^T \phi(Q)$$

Where  $\phi(x)^T \phi(y)$  is a kernel approximator and the scaling factor is omitted

- Linearized Attention greatly reduces computational cost
- There has been a lot of work in this field which has shown promise in recent years
- E.g. Retnet, Infini-Attention, Hedgehog etc.

Dual form of linear attention and gradient decent:

Let  $f_W(x) = Wx$  be a linear function. Given gradient descent with  $l_2$  loss,  $T$  training samples  $\{x_i, y_i\}_{i=1}^T$  and learning rate  $\eta$

$$W_1 x = \left( W_0 - \eta \nabla \frac{1}{T} \sum_{i=1}^T l_2(f_W(x_i), y_i) \Big|_{W=W_0} \right) x$$

$$= W_0 x + \text{LinAttn}\left(\frac{\eta}{T} E, X, x\right)$$

$X = [x_1; \dots; x_T]$  is the matrix of inputs

$E = Y - W_0 X$  is the error matrix where  $Y = [y_1, \dots, y_T]$

## Observations:

- In-context learning forms a type of meta-optimizer on the query resembling gradient descent with specific training data for linearized transformers
- Statement isn't constrained to specific regression settings and values for  $W_Q, W_K, W_V$
- $W_V p'_i$  is intuitively the “value” which we place upon token  $W_Q p'_i$ . Here we place emphasis on  $W_K p'_i$  rather than the query token itself.

## Extension to Multiple Layers:

- Consider a more realistic model architecture with  $L$  layers stacked upon each other
- $f(x) = (T_L + I) \circ \dots \circ (T_1 + I)(x)$
- For all  $i$  Layer  $T_i$  is either a FFN layer or a linear self-attention layer  $T_i = LSA_i$
- Corresponding weight matrices  $W_{Q_i}, W_{K_i}, W_{V_i}$
- $I$  is the identity function

### Algorithm 1: ICL imitation algorithm

```

1: input:  $f_1$  and  $[p'_m, \dots, p'_1, p_n, \dots, p_1]$ 
2: for  $i \in \{1, \dots, L\}$ 
  IF  $T_i$  is a FFN with residual connection, return
   $[p'_m, \dots, p'_1, p_n, \dots, p_1] = (T_i + I)([p'_m, \dots, p'_1, p_n, \dots, p_1])$ 

  ELSE  $T_i = LSA_i$ 
    (a) construct matrix  $W_0 = W_{base,i}([p_n, \dots, p_1])$ 
    (b) Update the linear functional  $f(x) = W_0 x$  with a single step of gradient descent with learning rate  $m$  and training set  $\{\phi(\cdot), W_{V_i} p'_j + W_0 \phi(W_{K_i} p'_j)\}_{j=1}^m$  such that the updated weights are  $W_1$ 
    (c)  $[p'_m, \dots, p'_1, p_n, \dots, p_1] = W_1 \phi(W_{Q_i} [p'_m, \dots, p'_1, p_n, \dots, p_1]) + [p'_m, \dots, p'_1, p_n, \dots, p_1]$ 
  
```

### Theorem 2:

For a model  $f_1$  described above and a prompt  $[p'_m, \dots, p'_1, p_n, \dots, p_1]$ , Algorithm 1 produces the same output as  $f_1(p'_m, \dots, p'_1, p_n, \dots, p_1)$

## Connection to greedy layer-wise algorithms

- Algorithm 1 is a recursive algorithm The labels are generated from the inputs themselves
- Algorithm 1 takes the form of a greedy unsupervised layer-wise pretraining algorithm (GLT)
- Real life phenomena shows degree of similarity:
  - GLT are observed to achieve quick convergence, ICL only applies a single step of gradient descent
  - GLT shown to help learn internal representations that represent higher level abstractions, ICL has similar properties
- This may motivate future steps forward with ICL:
  - This suggests that ICL may be a form of initialization which projects the model into a specific space
  - May want to have fixed ICL terms to guide fine-tuning steps, serving as the initialization. This actually begins to resemble instruction tuning
  - Motivates the construction of a unified theoretical framework combining fine-tuning methods with ICL and instruction tuning