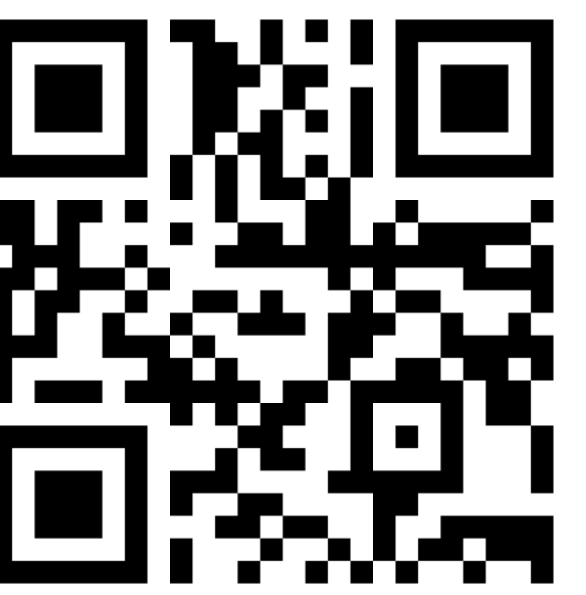
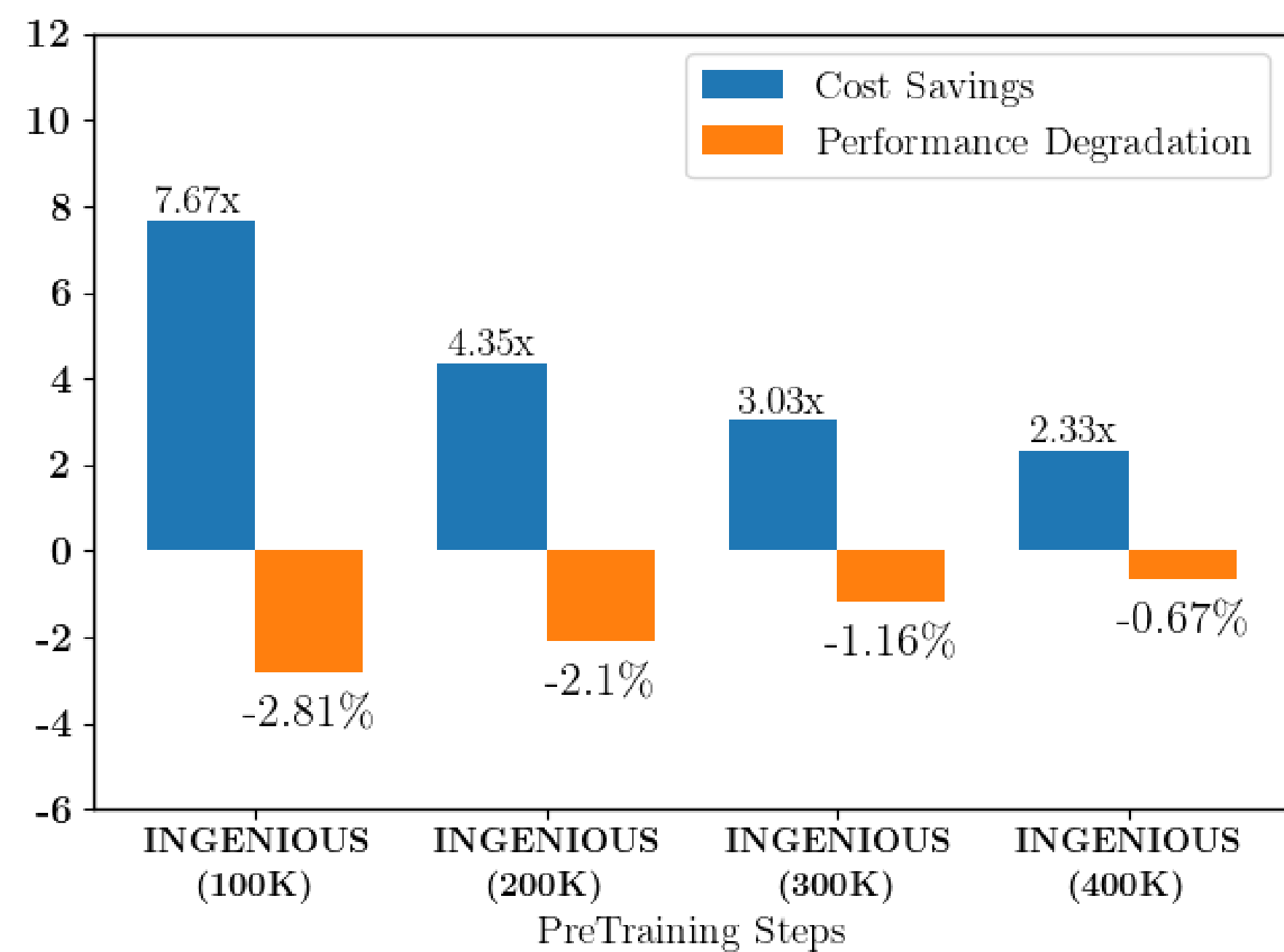


# INGENIOUS: Using Informative Data Subsets for Efficient Pre-Training of Language Models

H S V N S Kowndinya Renduchintala<sup>1</sup>, Krishnateja Killamsetty<sup>2</sup>, Sumit Bhatia<sup>1</sup>  
Milan Aggarwal<sup>1</sup>, Ganesh Ramakrishnan<sup>3</sup>, Rishabh Iyer<sup>2</sup>, Balaji Krishnamurthy<sup>1</sup>

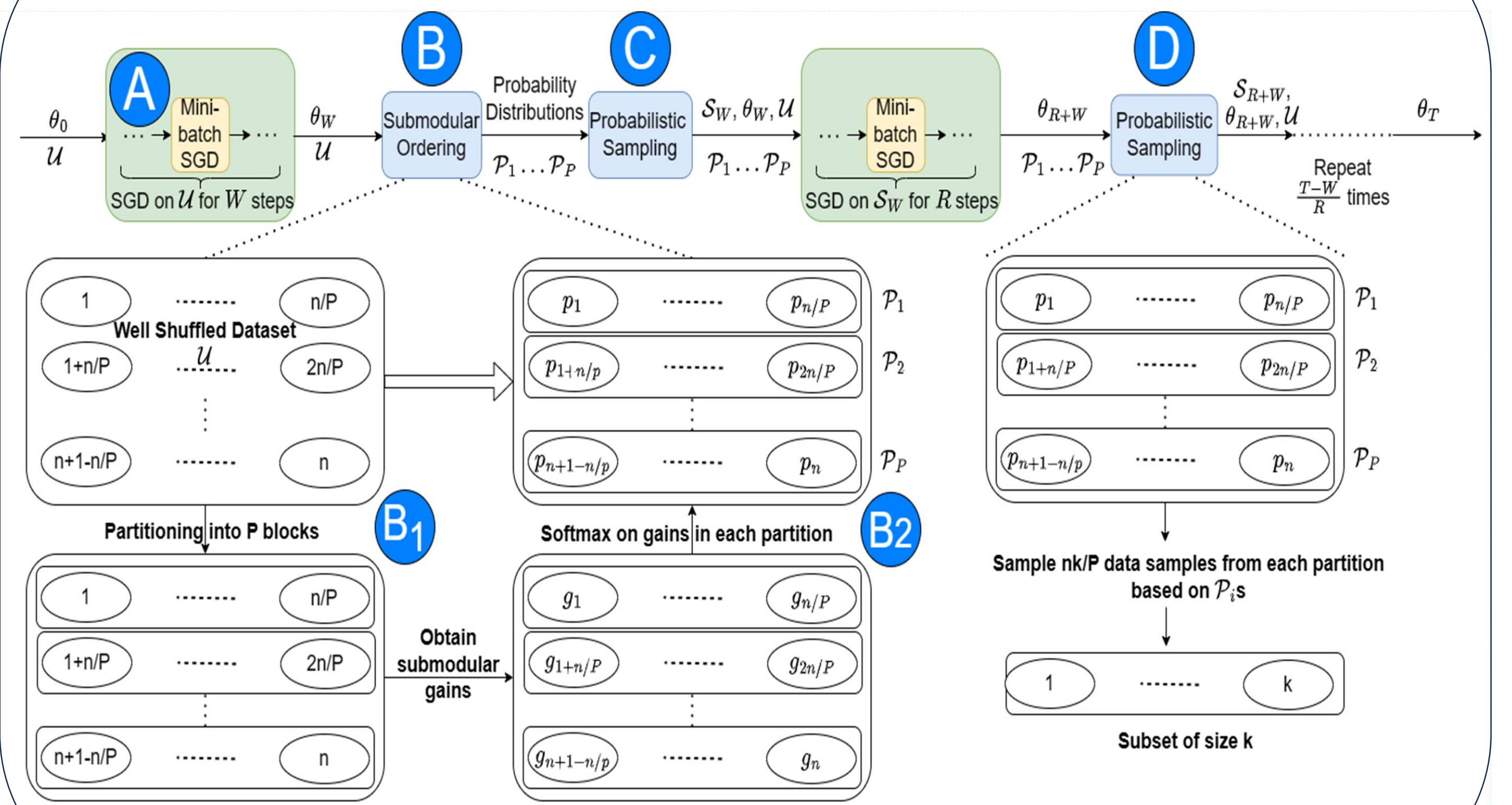


## Summary



INGENIOUS is an effective method to select informative subsets for efficient training of language models, based on submodular optimization.

## Framework



## Submodularity

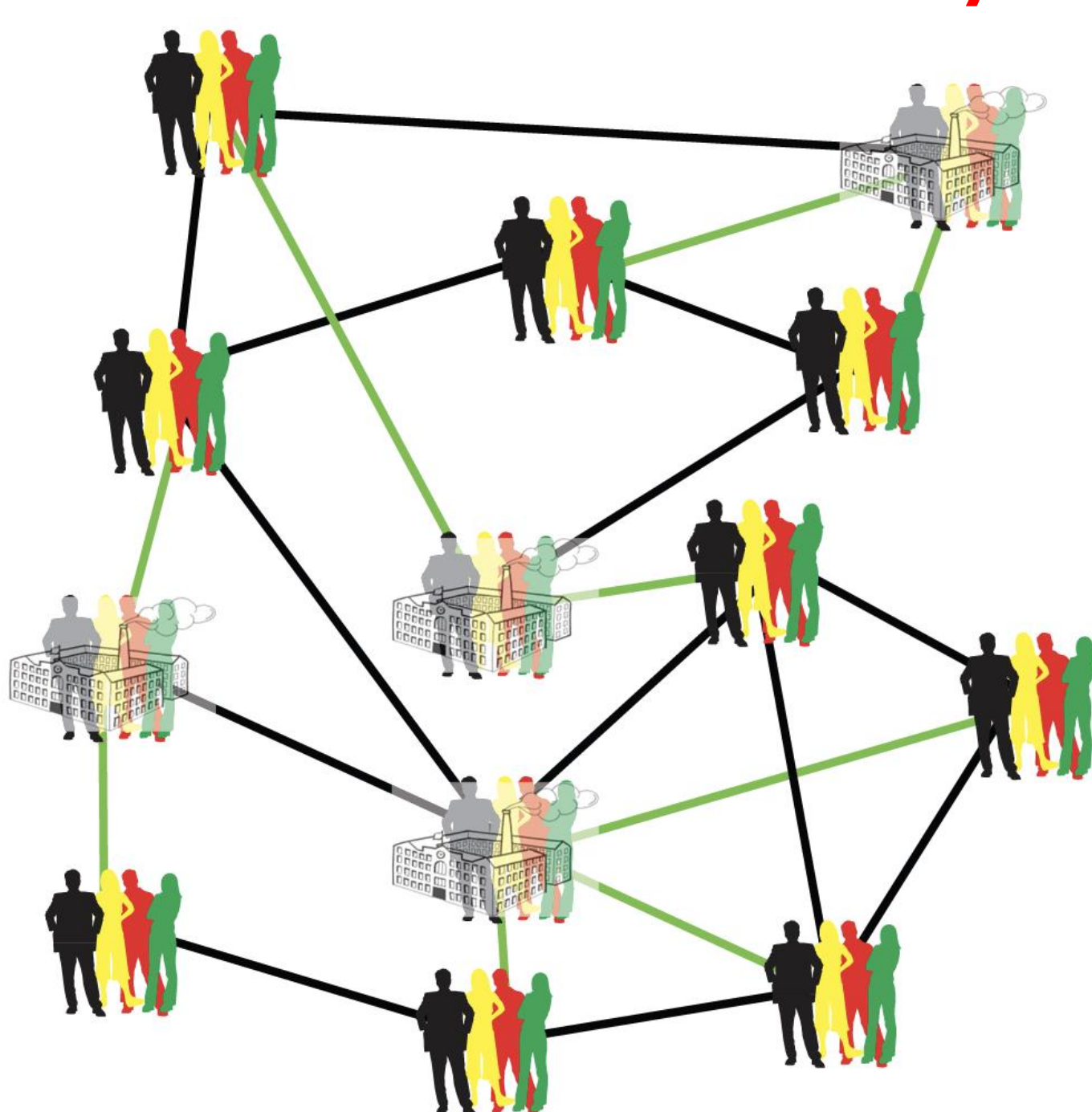
A set function  $f: 2^{\mathcal{V}} \rightarrow \mathbb{R}$  is called a submodular function if the following property is satisfied:

$$f(A \cup \{v\}) - f(A) \geq f(B \cup \{v\}) - f(B) \quad \forall A \subseteq B \subseteq \mathcal{V}; v \in \mathcal{V} \setminus B$$

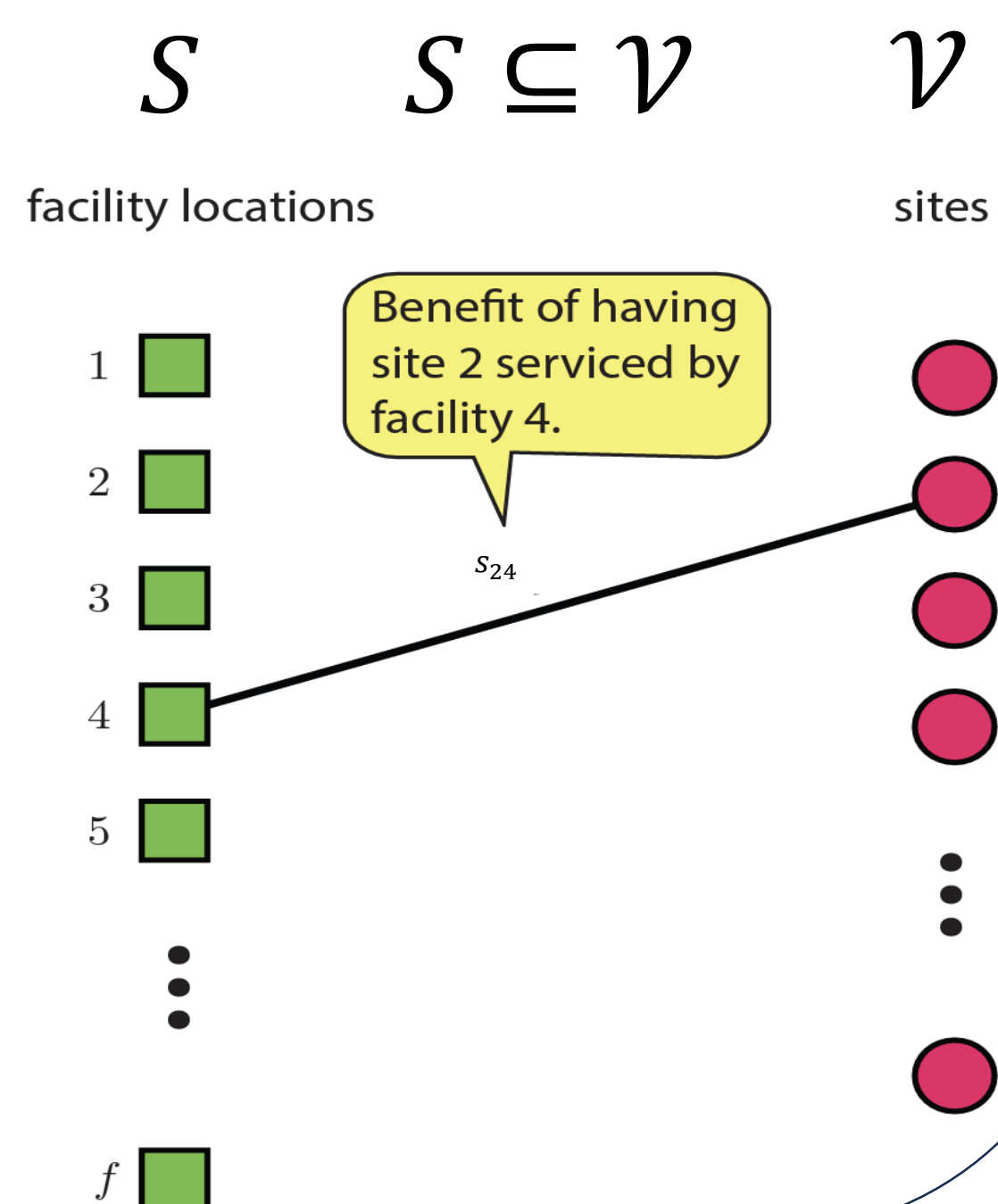
Example: If  $f$  denotes the consumer costs, **submodularity** expresses the following property of  $f$ :

$$f(\text{burger}) - f(\text{burger}) \geq f(\text{burger, fries}) - f(\text{burger})$$

## Facility Location

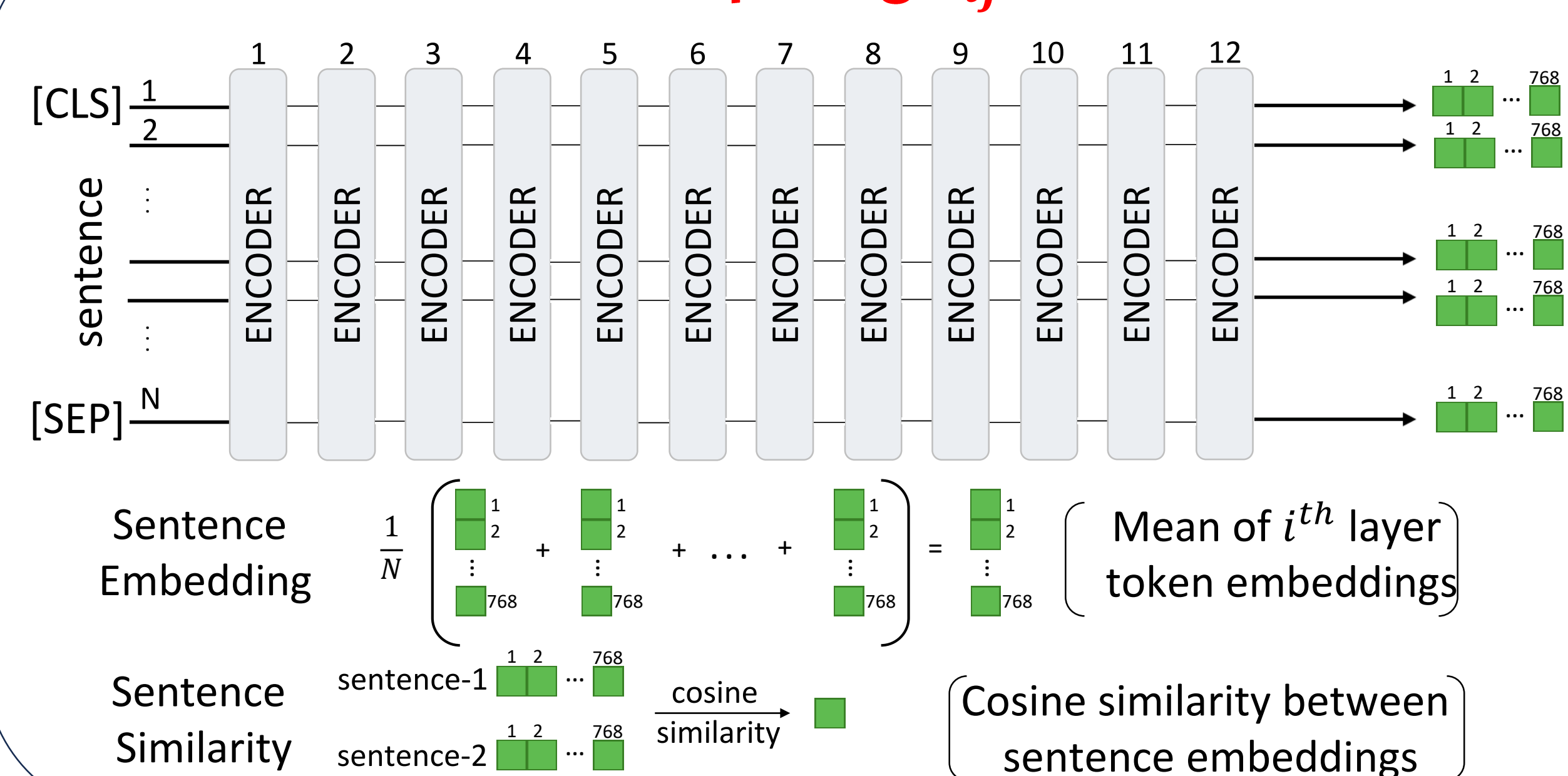


As efficiently as possible, place facilities at certain locations to satisfy sites having various demands.

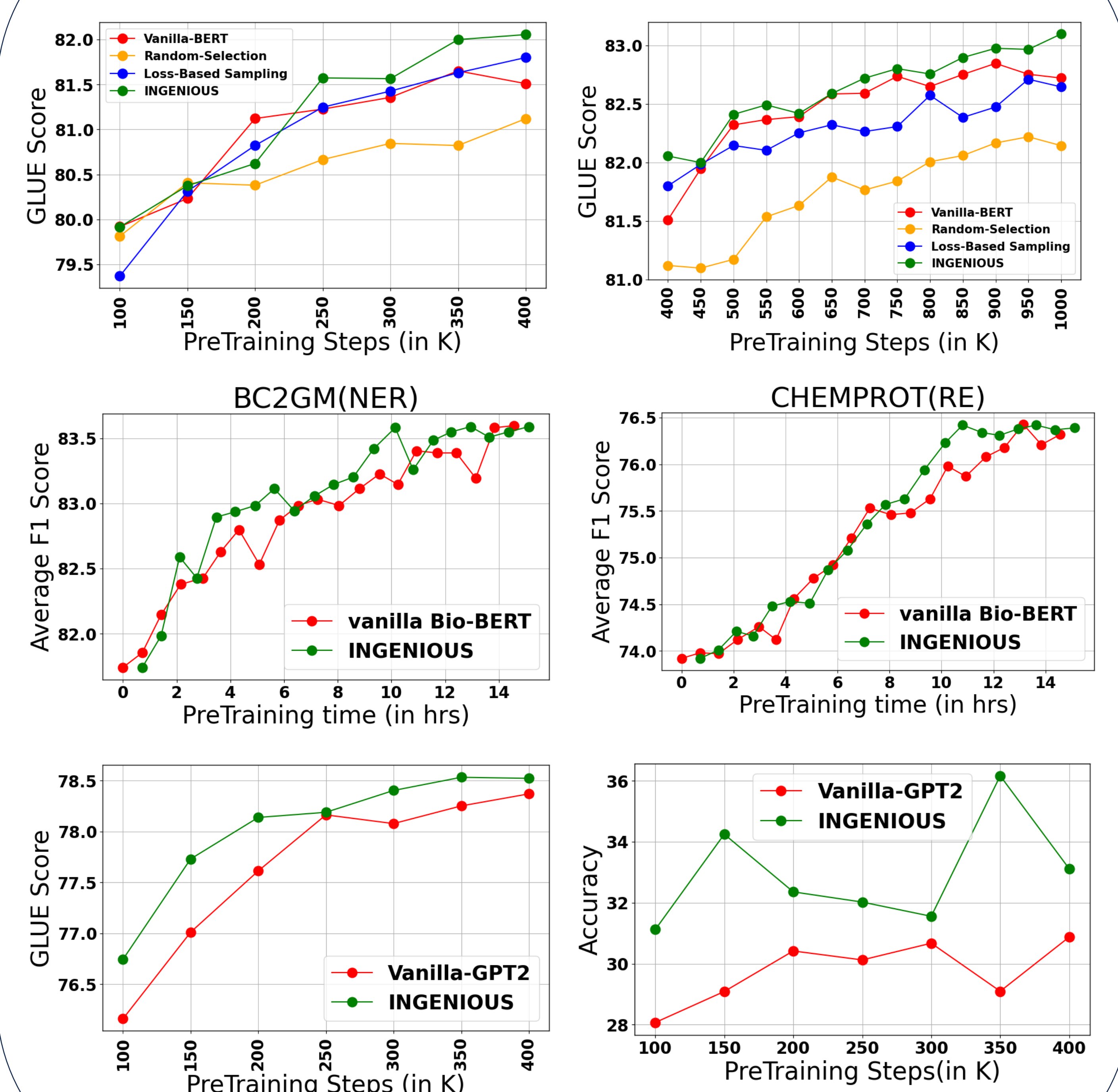


$$f_{FL}(S) = \sum_{i \in \mathcal{V}} \max_{j \in S} s_{ij}$$

## Computing $s_{ij}$ s



## Results



## Limitations & Future Directions

- Experiments pertaining to INGENIOUS framework are performed on relatively small language models compared to Llama or GPT-3. Future work could extend the framework to huge training corpora (~trillions of tokens) which are commonly used today.
- INGENIOUS can be extended to multi-modal settings where images and/or knowledge graphs can be brought in.
- Submodular measures such as Mutual Information can be used to efficiently train domain-specific language models.

<sup>1</sup>Adobe Inc., India

<sup>2</sup>University of Texas at Dallas, USA

<sup>3</sup>Indian Institute of Technology Bombay, India