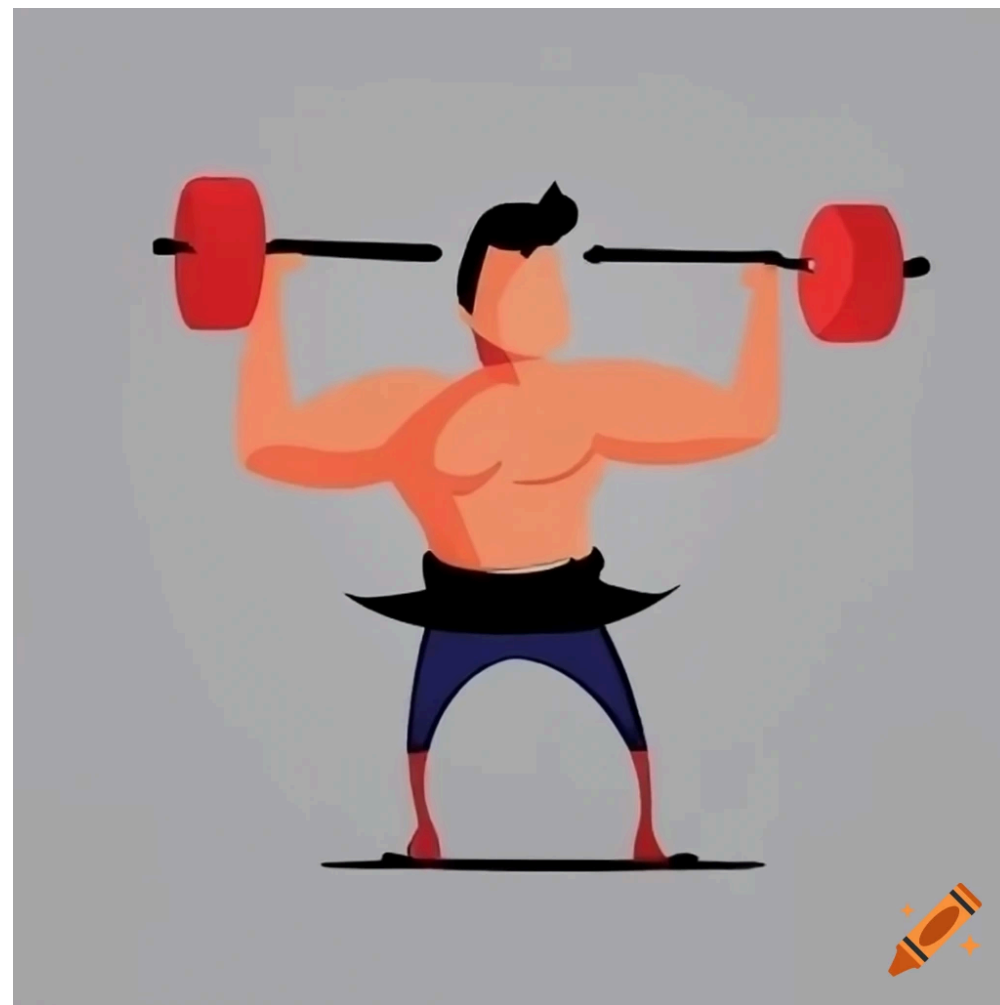# What is LIFT

Language-Interfaced Fine-Tuning for Non-Language Machine Learning Tasks

# 1. Background:

Pre-trained transformer language models have shown great success in natural language processing and protein/molecule design. This study explores their application in generating inorganic material compositions, a novel approach with potential for high-throughput materials discovery.

# 2. Challenge:

Navigating the vast design space of inorganic materials for new discoveries is costly and challenging. Pre-trained transformer models are leveraged to address the complexity of material composition generation, aiming for chemically valid and structurally stable results.

3. Current Status:

The study demonstrates the effectiveness of transformer models in generating hypothetical material compositions. Training different transformer models on diverse datasets yields high validity percentages, uniqueness, and potential for new materials discovery, opening new possibilities for materials design.

4. Method:

Seven transformer language models, including GPT, GPT-2, and RoBERTa, are trained on various material datasets. Validity, uniqueness, recovery rate, and novelty are evaluated to assess generative performance. Hyperparameters are tuned, and formation energies are predicted using DFT calculations.

5. Result:
   Transformer-based models exhibit high validity and novelty in generating material compositions. MT-GPTJ shows the best performance, closely followed by MT-GPT2 and MT-GPTNeo. MT-BART and MT-RoBERTa lag slightly behind. Training set size and dataset quality significantly impact model performance.

6. Conclusion:

Transformer language models show promise in generative materials design, capable of learning chemical patterns from training datasets. Despite biases and challenges in generating specific compositions, these models offer a new approach for materials exploration and potential discovery.

7. Outlook:

Future research may focus on addressing biases in composition generation, optimizing model training on diverse datasets, and integrating crystal structure predictions for synthesized materials. Continued development of transformer models for materials design could lead to significant advancements in high-throughput materials discovery.

# Thank you!