

EDAN70 Project: NLP Relation Extraction in Biomedical Literature

Jacob Krucinski

Supervisor: Sonja Aits

Co-supervisor: Rafsan Ahmed

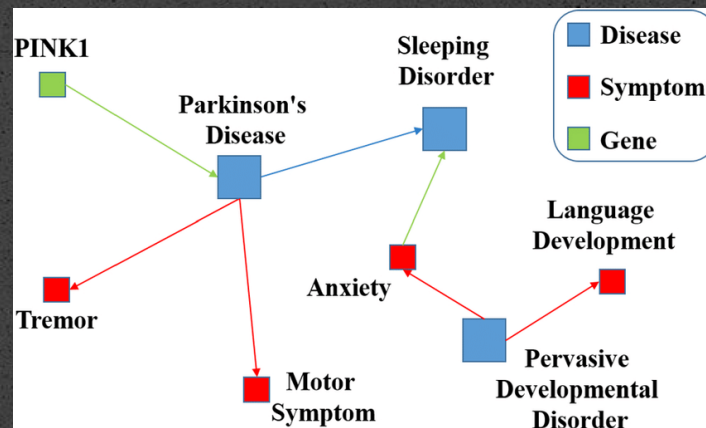
Cell Death, Lysosomes and Artificial Intelligence Group, Lund University

<https://github.com/Aitslab/BioNLP>

ja6750kr-s@student.lu.se

GitHub: [@savage-hacker14](https://github.com/savage-hacker14)

LinkedIn: <https://www.linkedin.com/in/jkrucinski/>

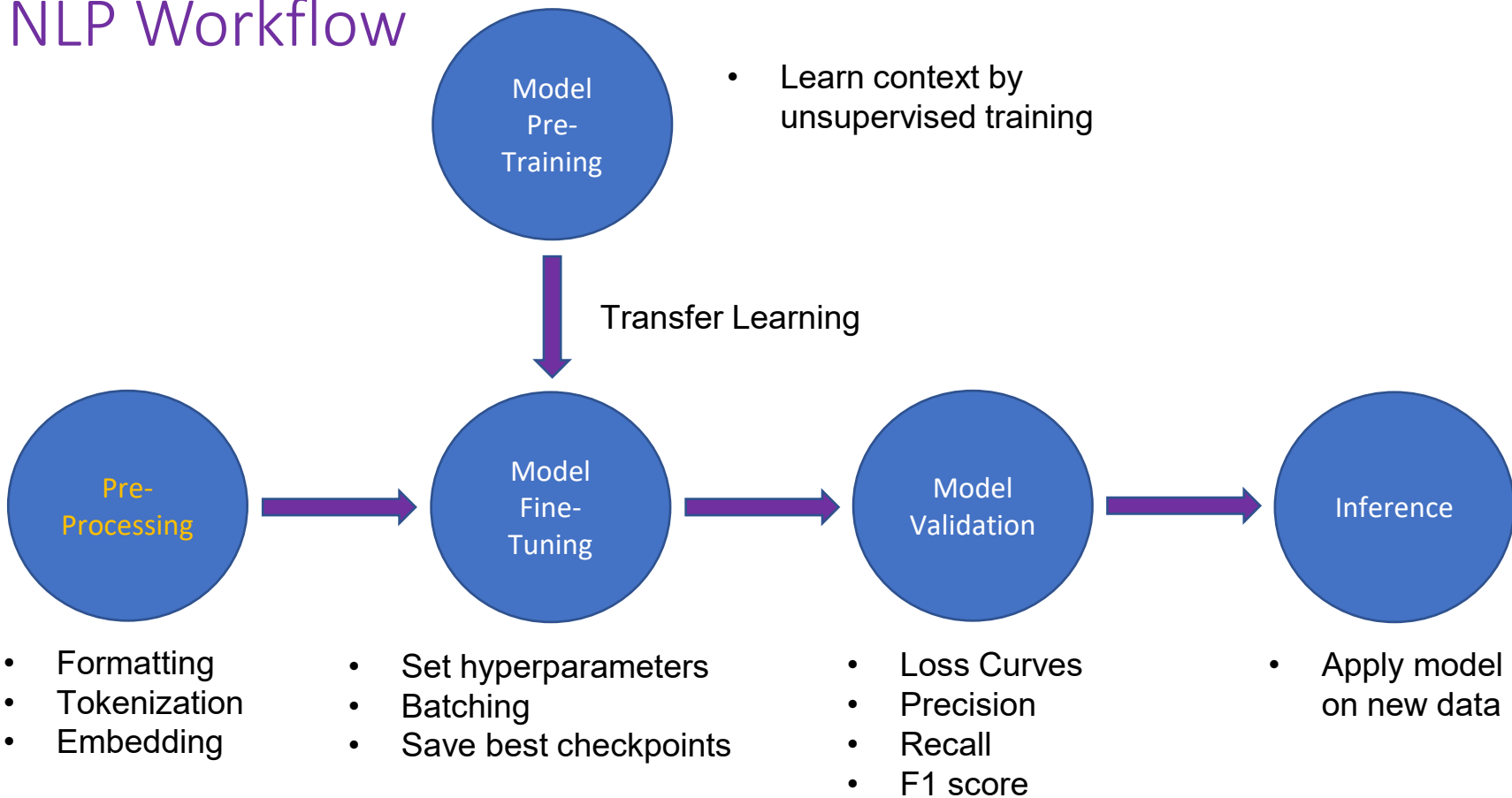


Introduction to NLP

Biomedical Domain



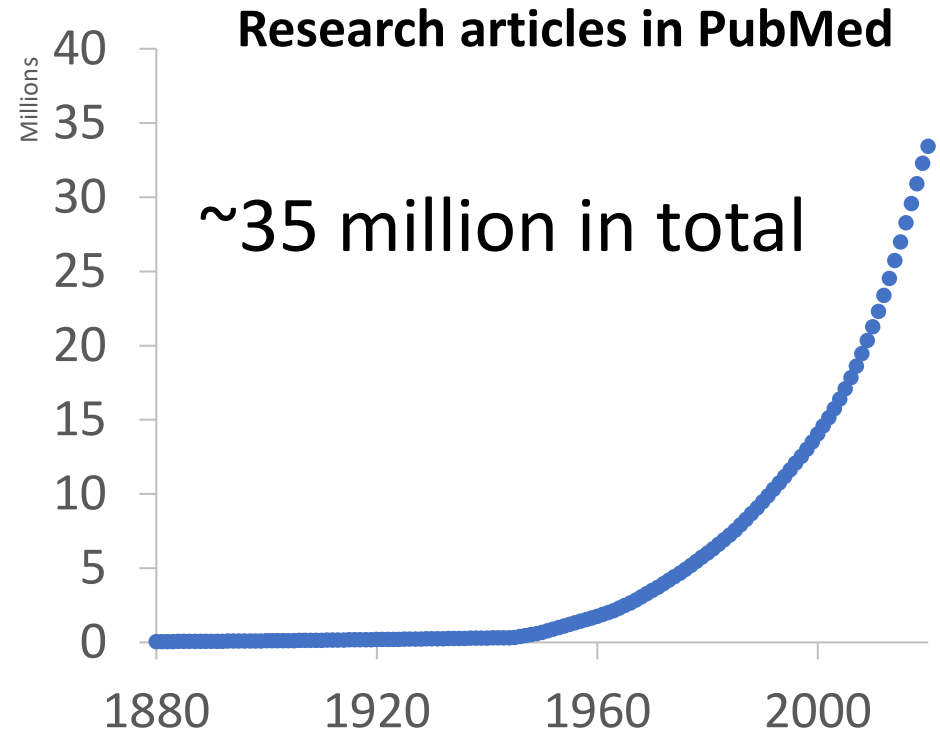
NLP Workflow





Humans can no longer process the accumulated medical knowledge

Importance in the Biomedical Domain



Models: BioGPT & SciBERT

GPT: Generative Pre-Trained Transformer
 BERT: Bidirectional Encoding Representations from Transformers

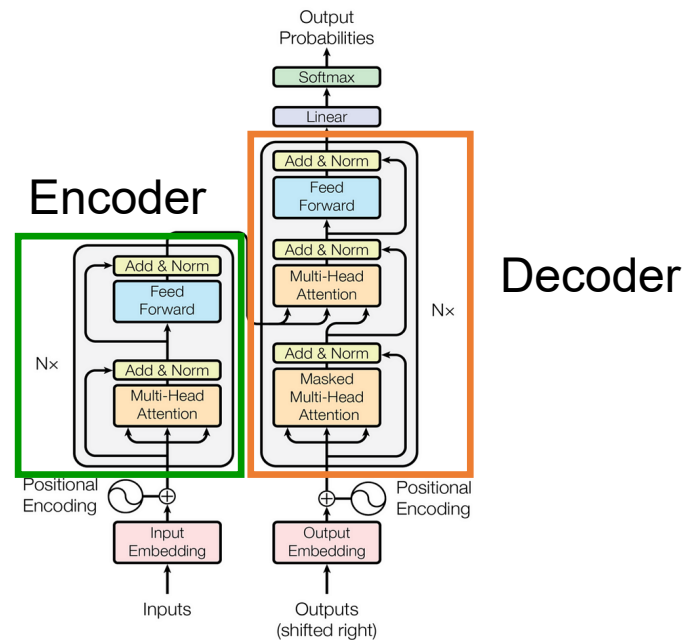
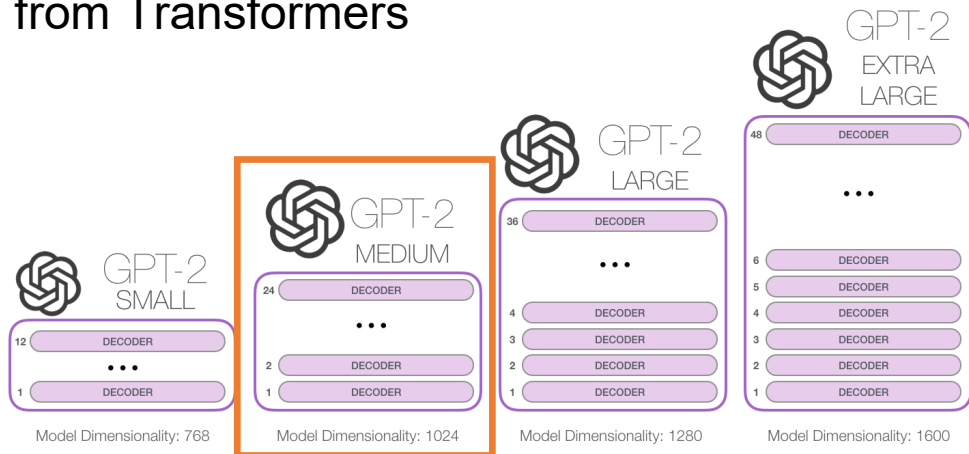


Figure 1: The Transformer - model architecture.

Tasks: End-to-end RE, Q&A, Document Classification, and text generation

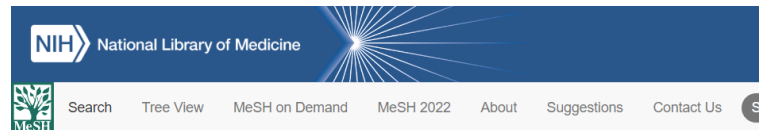
$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Sub-project 1

BioGPT Hyperparameter Tuning

Corpus: BioCreative 5 Chemical-Disease Relations

Task Dataset	Articles	Chemical		Disease		CID relation
		Mention	ID	Mention	ID	
Training	500	5,203	1,467	4,182	1,965	1,038
Development	500	5,347	1,507	4,244	1,865	1,012
Test	500	5,385	1,435	4,424	1,988	1,066



Famotidine MeSH Descriptor Data 2023

Details	Qualifiers	MeSH Tree Structures	Concepts
MeSH Heading	Famotidine		
Tree Number(s)	D02.886.675.215 D03.363.129.708.215		
Unique ID	D015738		
RDF Unique Identifier	http://id.nlm.nih.gov/mesh/D015738		
Scope Note	A competitive histamine H2-receptor antagonist. Its main pharm		
Entry Term(s)	Famotidine Hydrochloride		
	MK-208		
	Pepcid		
	YM-11170		

8701013|t|Famotidine-associated delirium. A series of six cases.

8701013|a|Famotidine is a histamine H2-receptor antagonist used in inpatient settings for prevention of stress ulcers and is showing increasing popularity because of its low cost. Although all of the currently available H2-receptor antagonists have shown the propensity to cause delirium, only two previously reported cases have been associated with famotidine.

(TRUNCATED)

8701013	0	10	Famotidine	Chemical	D015738
8701013	22	30	delirium	Disease	D003693

(TRUNCATED)

8701013	CID	D015738	D003693
---------	-----	---------	---------



Experiment Hyperparameters

Train 1: 50 epochs

Model:

Dropout 0.1

Optimizer:

Adam Betas: (0.9, 0.98)

Learning rate: 1e-5 (w/ inv. sqrt scheduler)

Train 2: 50 epochs

Model:

Dropout 0.2

Optimizer:

Adam Betas (0.9, 0.99)

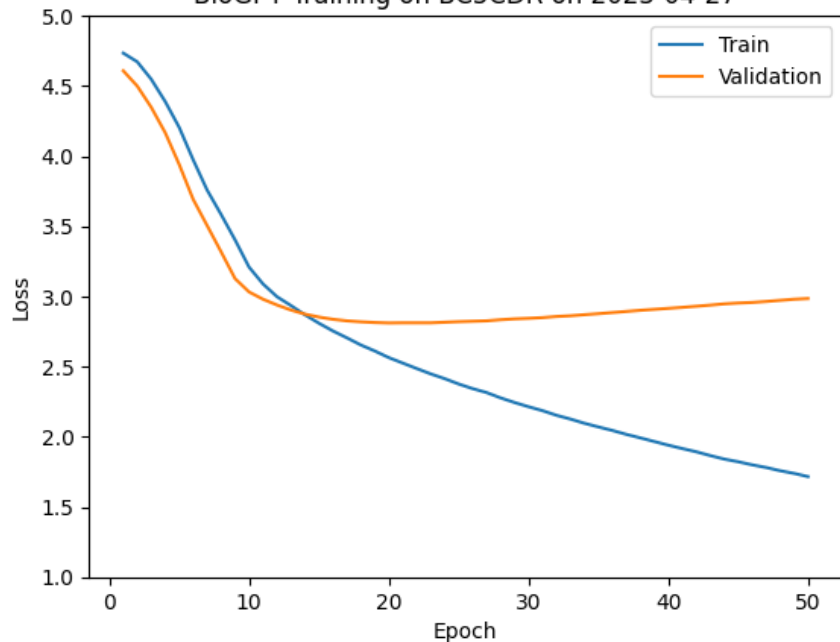
Learning rate: 2e-5 (same scheduler)



Train & Validation Loss Curves

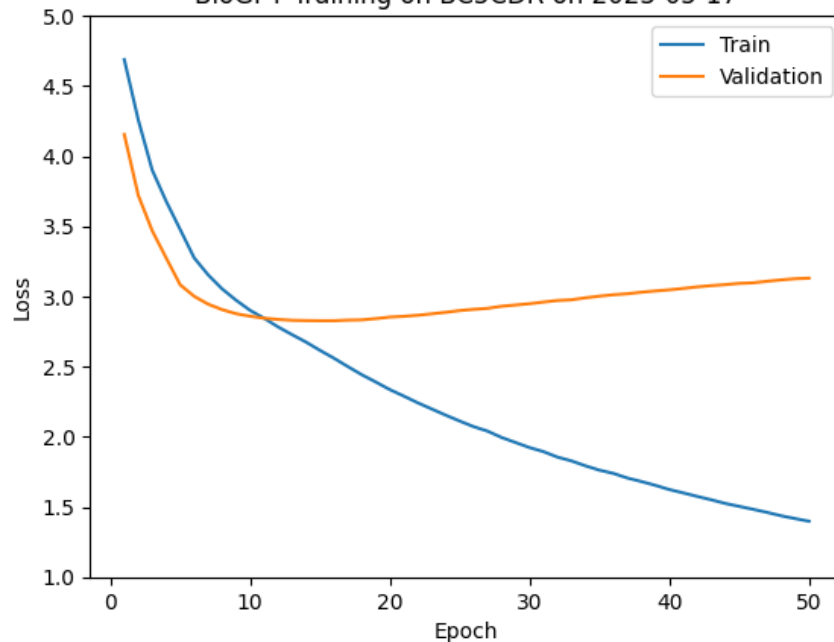
Experiment 1

BioGPT Training on BC5CDR on 2023-04-27



Experiment 2

BioGPT Training on BC5CDR on 2023-05-17





Performance Comparison to Original Paper

Metric	Original Paper	Experiment 1 (+/- change from paper)	Experiment 2 (+/- change from paper)
Precision	0.4944	0.4872 (-0.0072)	0.5176 (+0.0232)
Recall	0.4128	0.4118 (-0.0010)	0.4005 (-0.0123)
F1	0.4498	0.4464 (-0.0034)	0.4516 (+0.0018)

** using last 5 checkpoint averaging

Sub-project 2

BioGPT Free-form Text Relation Extraction



Inference Methods: Forward Pass & Text Generation

True Inference:

- Pre-processing: Tokenization, Byte Pair Encoding (BPE)
- Forward pass through the model
- Post-processing: DeBPE, Detokenization, etc.

Text Generation:

- a) Prompt engineering for NER
- b) “rel-is” prompting via manual NER

“rel-is” example: The relation between **dextropropoxyphene** and **mu-type opioid receptor** is **inhibitor**

Results: True Inference

Source	Free-form Text and Output	True Entities	True Relations
PubTator ID 439781	<p>IN: After a single oral dose of 4 mg/kg indomethacin (IDM) to sodium and volume depleted rats plasma renin activity (PRA) and systolic blood pressure fell significantly within four hours, suggesting hypotension.</p> <p>OUT (with my model): 227508 CID -1 -1 1.0</p> <p>OUT (with paper model): 227508 CID D012964 -1 1.0 (sodium)</p>	<p>C: Indomethacin D: hypotension</p>	<p>Indomethacin/ hypotension</p>
Custom	<p>IN: A patient perscribed 100 mg/ml cyclosporine suffered from high blood pressure, negative thoughts, and depression.</p> <p>OUT (same for my model and paper model): 227508 CID D016572 D006973 1.0 (cyclosporine <--> hypertension)</p>	<p>C: Cyclosporin D: Depression</p>	<p>Cyclosporin/ depression</p>



Results: Text Generation with “rel-is” Prompting

Source	Free-form Text and Output	True Entities	True Relations
Custom	<p>IN: A patient perscribed 100 mg/ml cyclosporine suffered from high blood pressure, negative thoughts, and depression.</p> <p>OUT: The relation between cyclosporine and depression is complex and may be influenced by several factors, including the dose of cyclosporine, the duration of cyclosporine therapy, and the presence of concomitant medications, such as lithium, which may affect the pharmacokinetics of cyclosporine and the development of cyclosporine-induced depression. (ABSTRACT TRUNCATED AT 250 WORDS) Cyclosporine-induced depression is a serious and potentially life-threatening side effect of cyclosporine therapy, and it is important to identify patients at risk for the development of cyclosporine-induced depression so that appropriate interventions can be instituted.</p>	<p>C: Cyclosporin D: Depression</p>	<p>Cyclosporin/depression</p>

Sub-project 3

SciBERT Fine-tuning with Combined
ChemProt & DrugProt Corpora



Corpus: ChemProt & DrugProt Combined

ChemProt: Chemical-Protein Interactions

- 1,820 PubMed abstracts, 16,075 total interactions

DrugProt: Drug-Protein Interaction (“Sister Corpus” for ChemProt)

- 4,250 abstracts, 21,035 train/dev interactions
- No CPR codes, but same relation groups

Train Statistics

Class	Count	Percentage
REGULATOR-NEGATIVE	10207	43.04%
INTERACTOR	8435	35.58%
REGULATOR-POSITIVE	3634	15.33%
PART-OF	1193	5.03%
NOT*	241	1.02%

Development Statistics

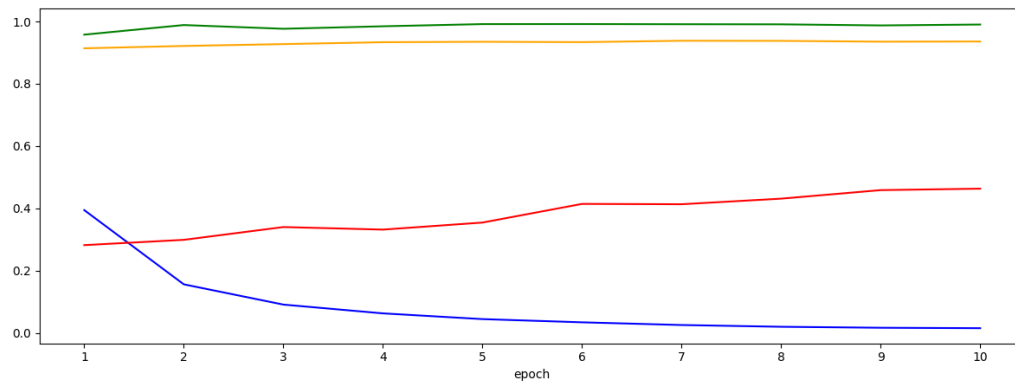
Class	Count	Percentage
REGULATOR-NEGATIVE	3004	41.06%
INTERACTOR	2594	35.45%
REGULATOR-POSITIVE	1134	15.50%
PART-OF	410	5.60%
NOT*	175	2.39%



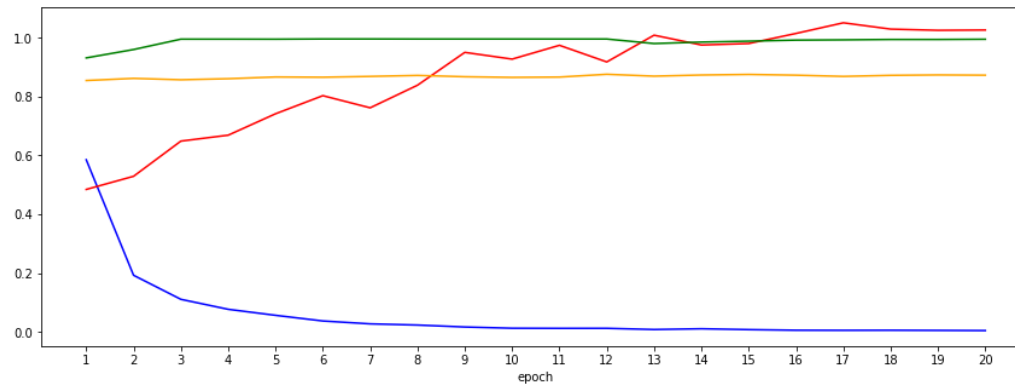
Results: Loss Curves



Combined Model



Baseline Model





Results: Precision, Recall, F1-Score

	Combined Model			Baseline Model	
Metric	Train score	Dev score	Dev change	Train score	Dev score
Precision					
macro	0.9925	0.9113	+0.0266	0.9955	0.8847
micro	0.9924	0.9385	+0.0577	0.9958	0.8808
weighted	0.9924	0.9382	+0.0584	0.9958	0.8798
Recall					
macro	0.9943	0.8822	+0.0252	0.9973	0.8570
micro	0.9924	0.9385	+0.0577	0.9958	0.8808
weighted	0.9924	0.9385	+0.0577	0.9958	0.8808
F1-score					
macro	0.9928	0.8902	+0.0292	0.9952	0.8610
micro	0.9924	0.9385	+0.0577	0.9958	0.8808
weighted	0.9924	0.9382	+0.0583	0.9958	0.8799

Conclusions & Key Takeaways

BioGPT Hyperparameter Tuning:

- Reduced training time to lowest validation loss from 22 to 14 epochs
- 2.3% increase in precision in 2nd experiment

BioGPT Free-form Relation Extraction:

- End-to-end RE can fail to extract entities and multiple relations
- NER prompting does *not* work well
- “Rel-is” prompting provides pathways/mechanisms for interactions

SciBERT Fine-tuning with Combined Corpus:

- The combined ChemProt/DrugProt corpus increased metrics by ~5%

Q&A