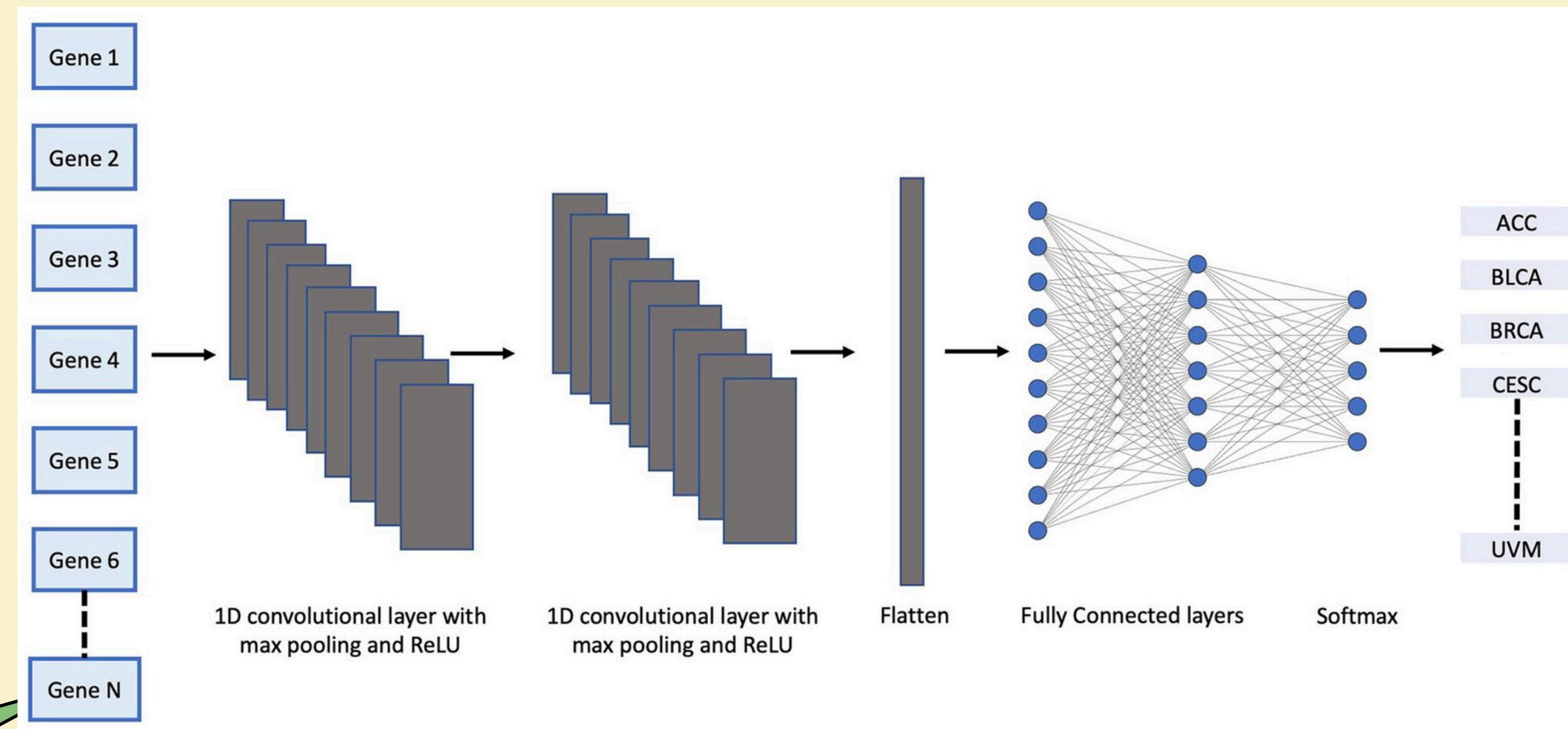


TULIP

Giorgi Kuchava

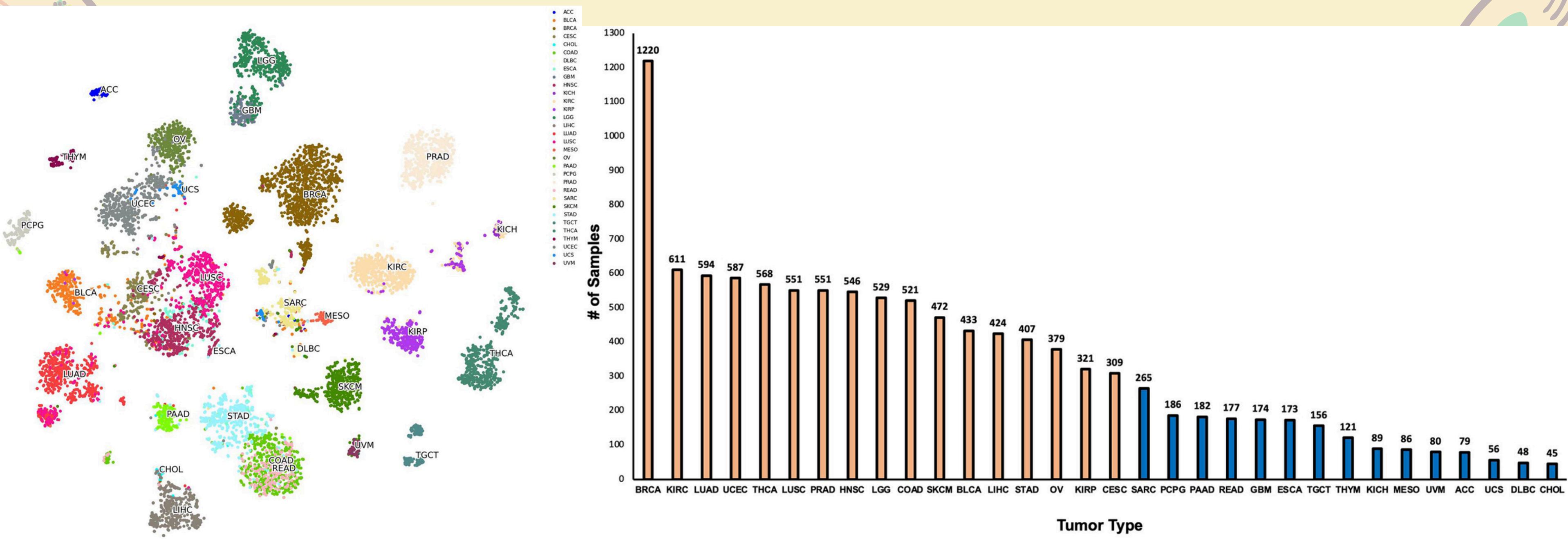
Giorgi Bendianishvili

TULIP: An RNA-seq-based Primary Tumor Type Prediction Tool Using Convolutional Neural Networks

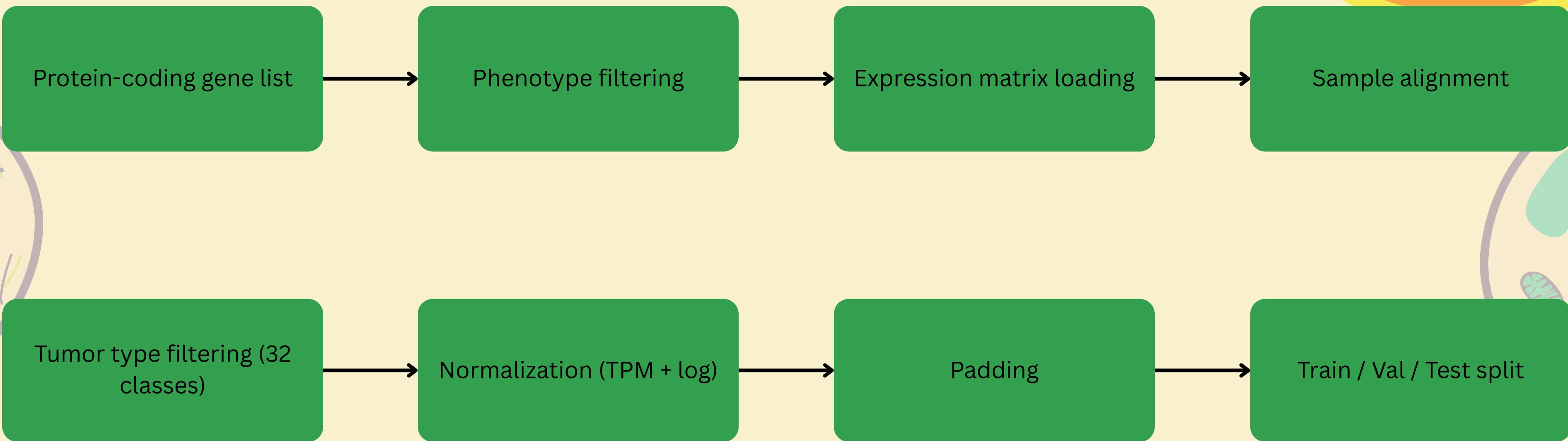


Tumor Types

19,758 protein coding genes as input and 32 primary tumor types



Preprocessing



Raw Data Structure

Gene ID	Sample1	Sample2	Sample3	Sample4	Sample5
ENSG1	32	51	0	73	11
ENSG2	104	82	63	22	5
ENSG3	0	0	44	55	61
ENSG4	91	13	22	0	33
ENSG5	77	88	99	11	22

Primary Tumor Filtering

Sample	Type
S1	Primary Tumor
S2	Normal
S3	Metastatic
S4	Primary Tumor
S5	Normal

Sample	Type
S1	Primary Tumor
S4	Primary Tumor

Tumor-Type Filtering

Sample	Cancer Type
S1	Breast
S2	Lung
S3	Rare Sarcoma
S4	Thyroid
S5	Unknown

Sample	Cancer Type
S1	Breast
S2	Lung
S4	Thyroid

Protein-Coding Gene Filtering(chunking)

Gene	Sample1	Sample2
ENSG1 (coding)	32	51
ENSG2 (non-coding)	104	82
ENSG3 (coding)	0	44
ENSG4 (lncRNA)	91	13
ENSG5 (coding)	77	88

Gene	Sample1	Sample2
ENSG1	32	51
ENSG3	0	44
ENSG5	77	88

FPKM → TPM → log10

Gene	Value	Gene	TPM	Gene	log10(TPM)
G1	50	G1	500 000	G1	5,7
G2	30	G2	300 000	G2	5,48
G3	20	G3	200 000	G3	5,3

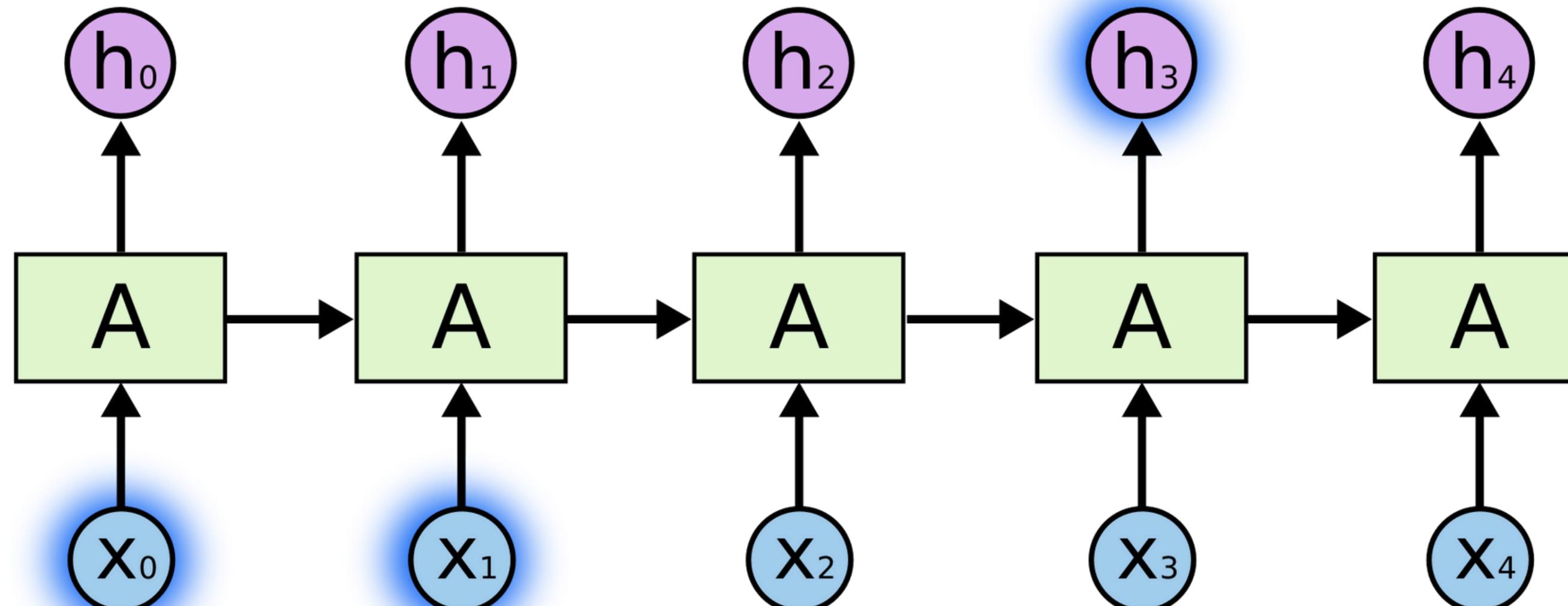
Train/Val/Test split

80% Train
10% Validation
10% Test

LSTM



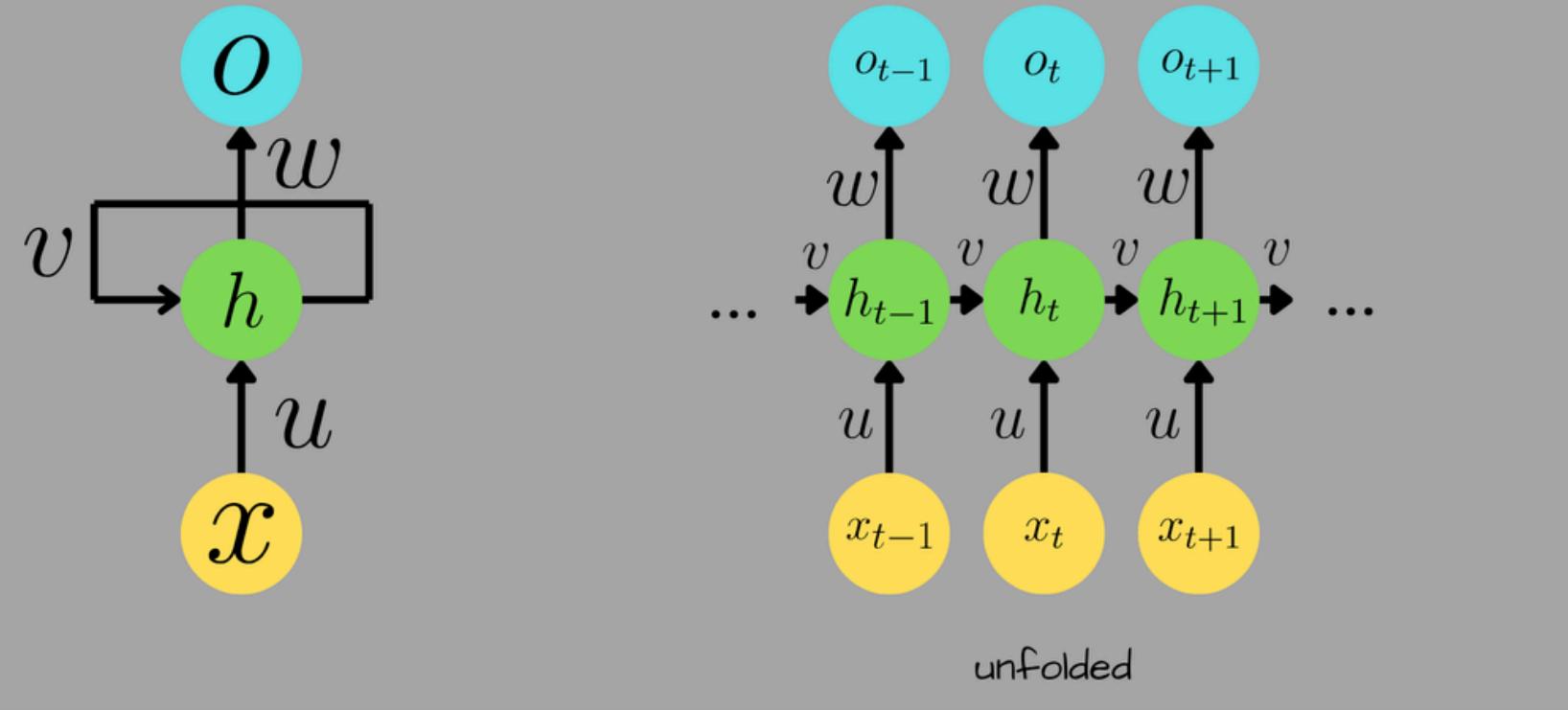
Main Problem with RNN



For example: “The movie that I watched yesterday was amazing”

Problems with RNN

RNN - architecture

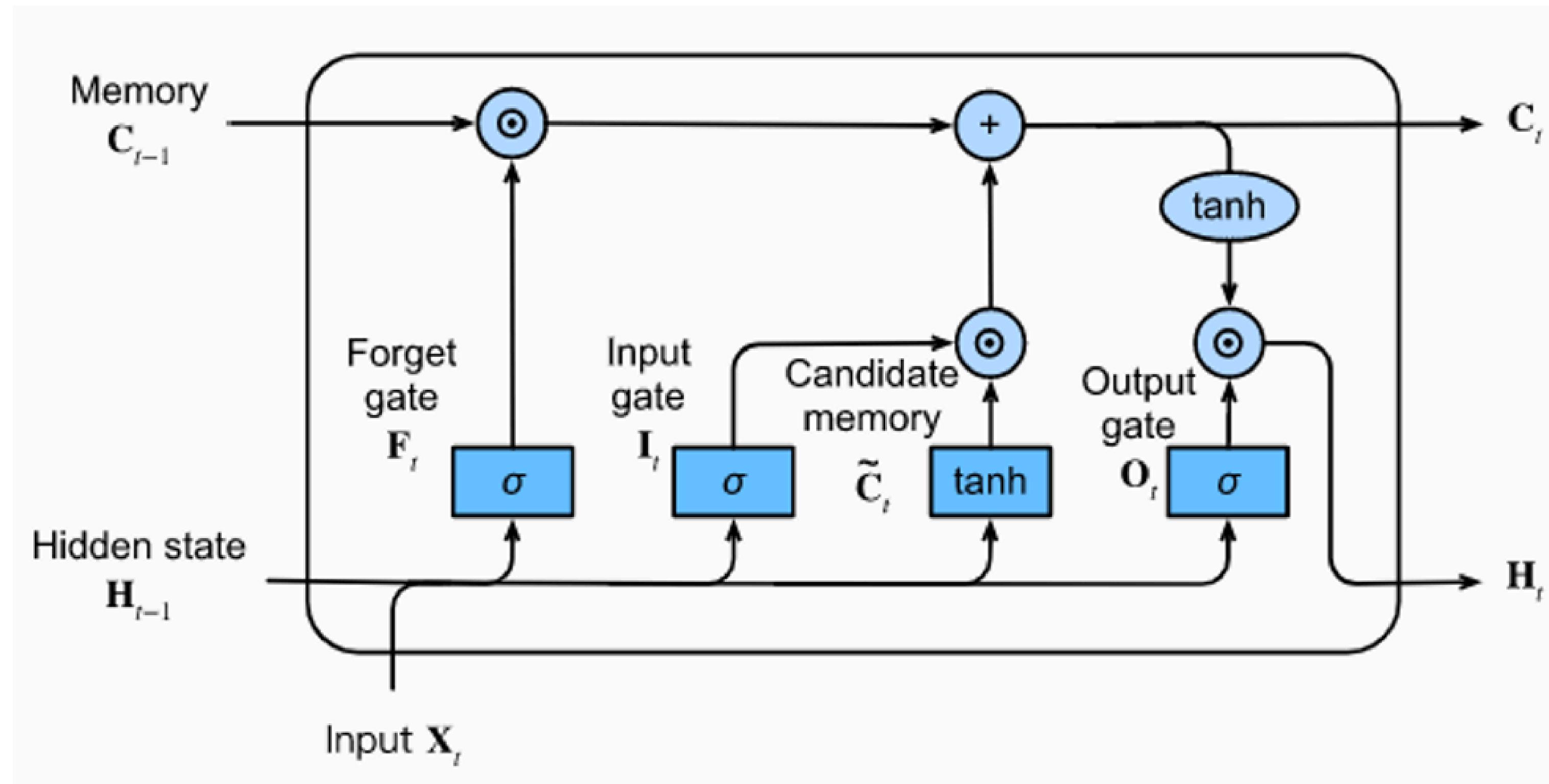


The h gets overwritten
every time

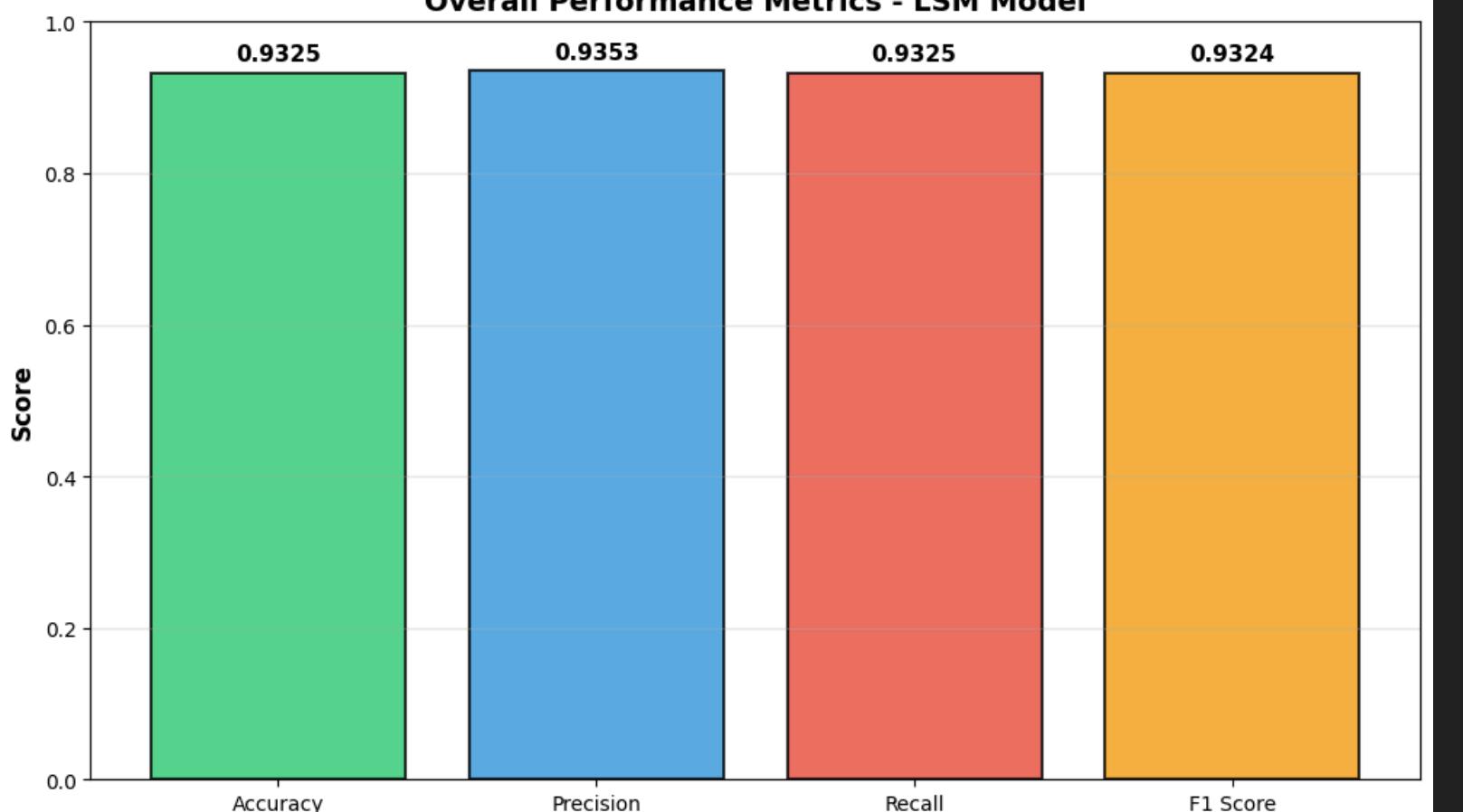
Exploding Gradients if
weights > 1

Vanishing Gradients if
weights < 1

LSTM

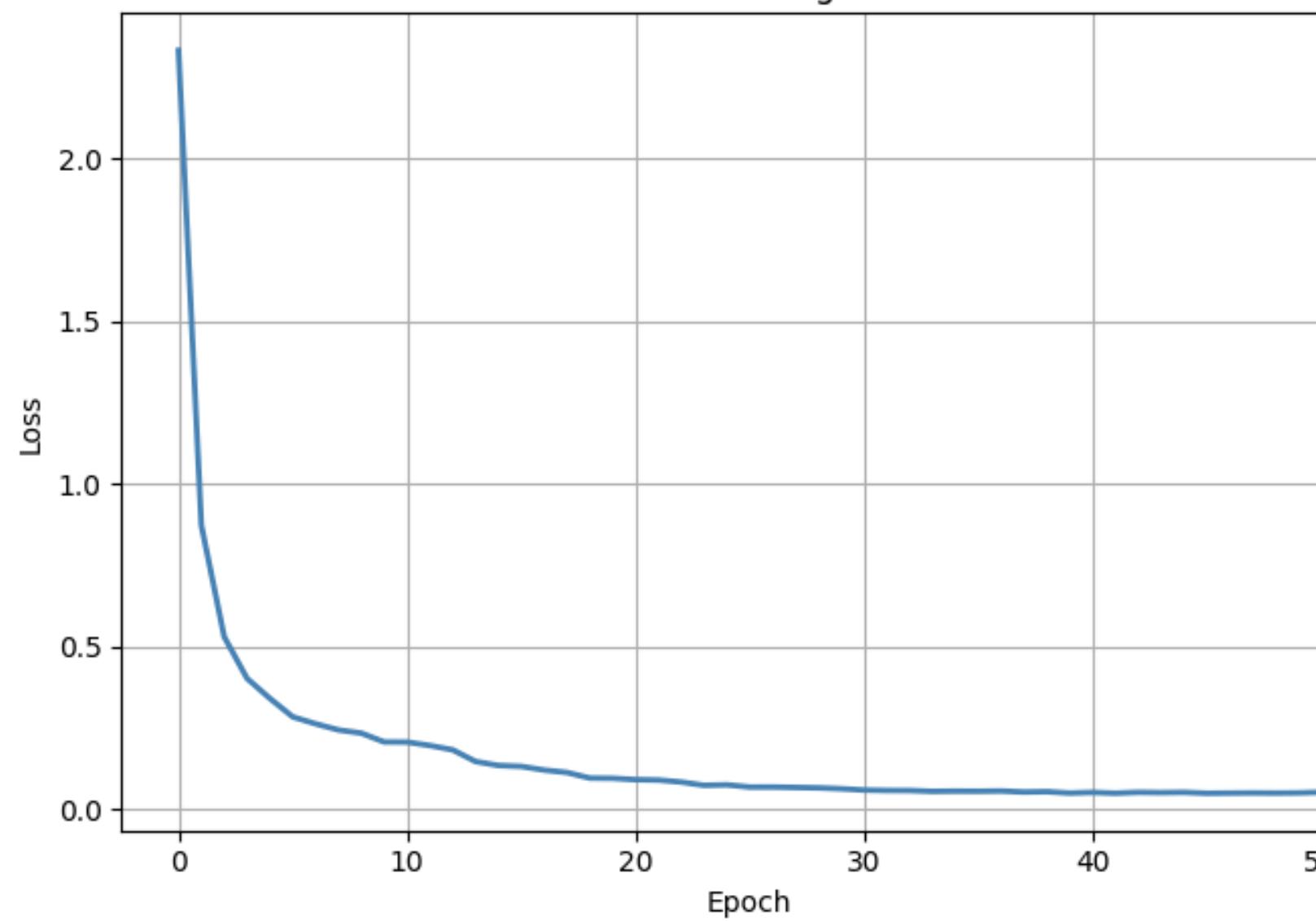


Overall Performance Metrics - LSM Model

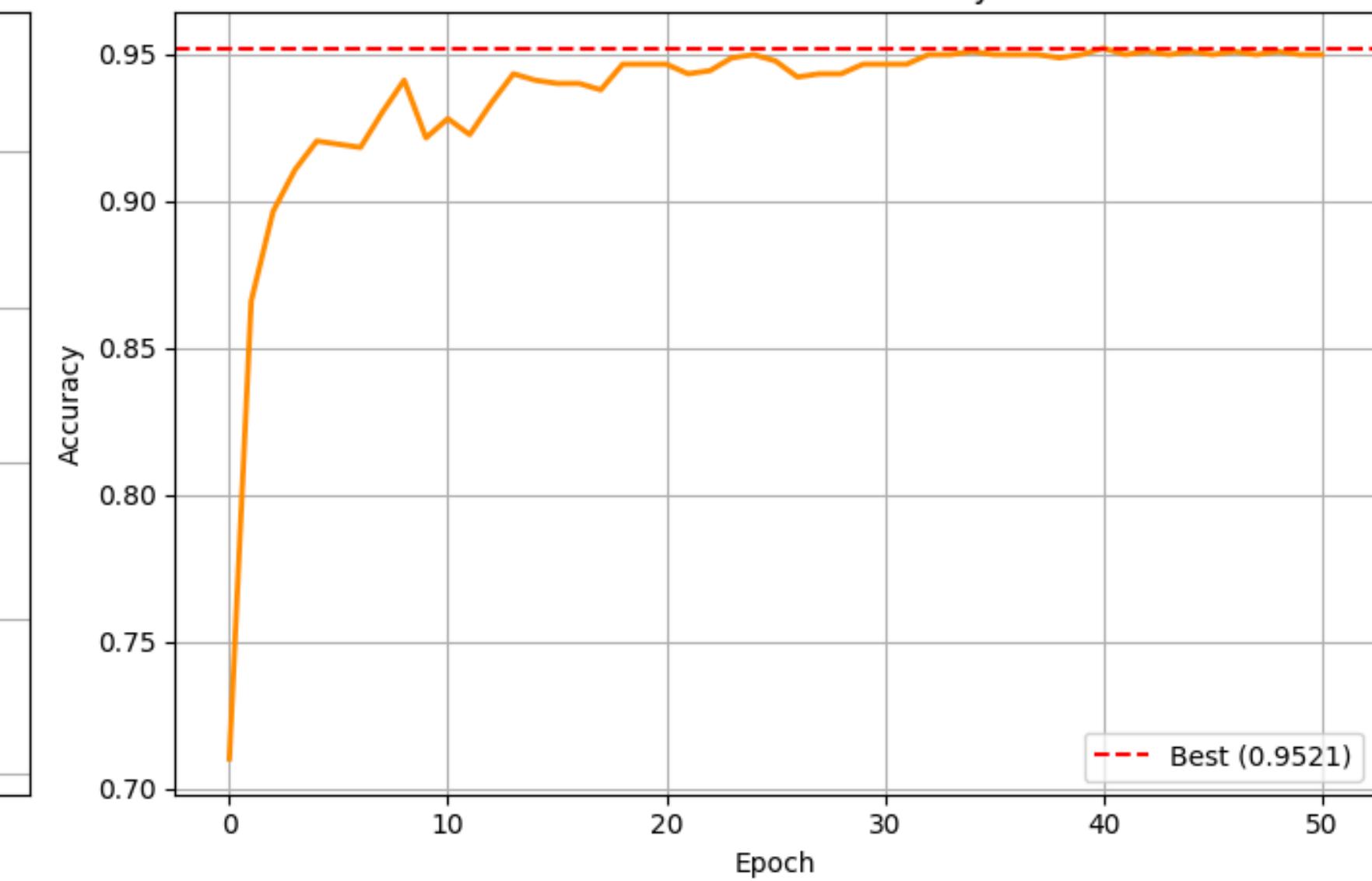


Training

LSM — Training Loss

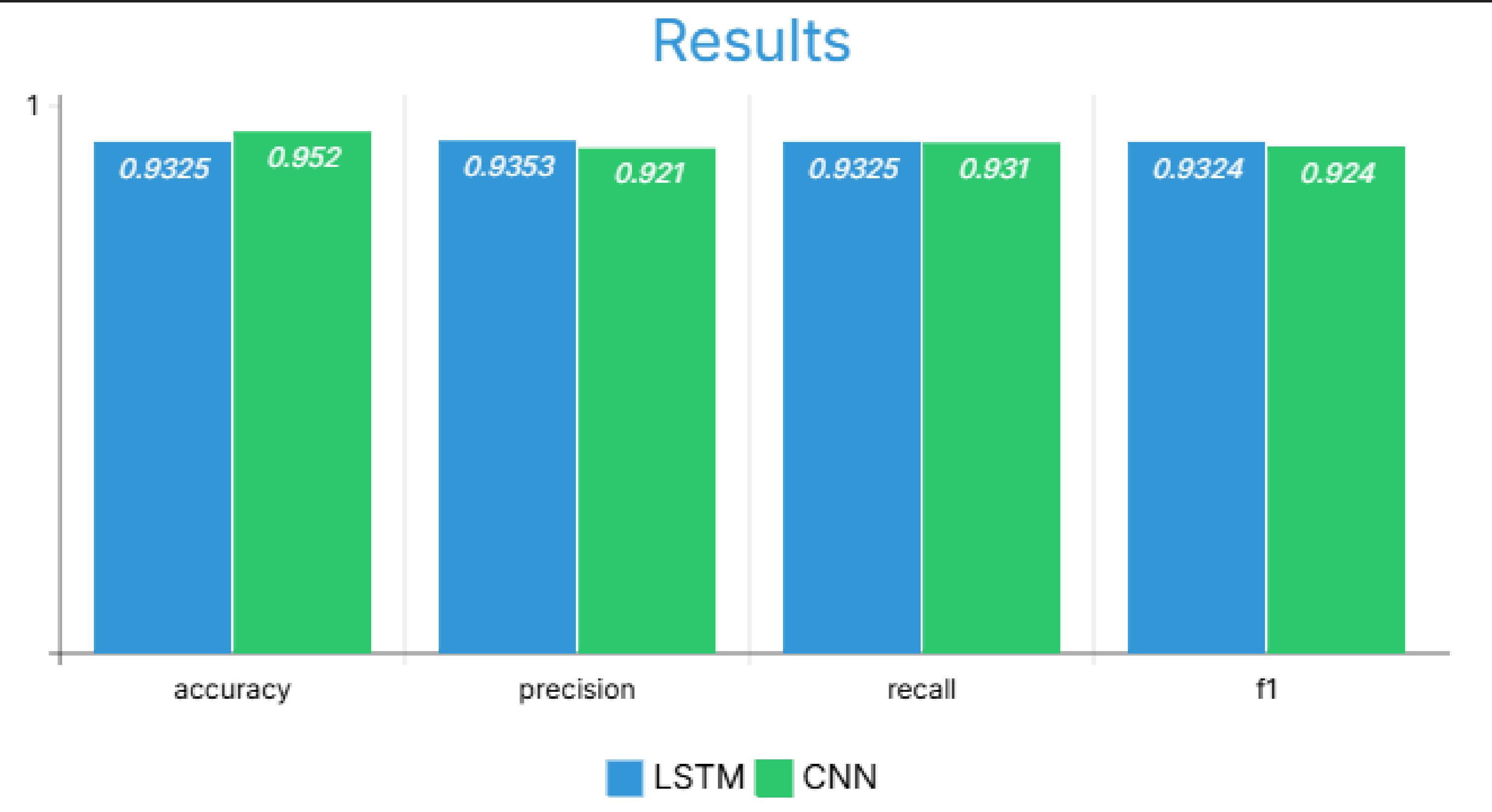


LSM — Validation Accuracy

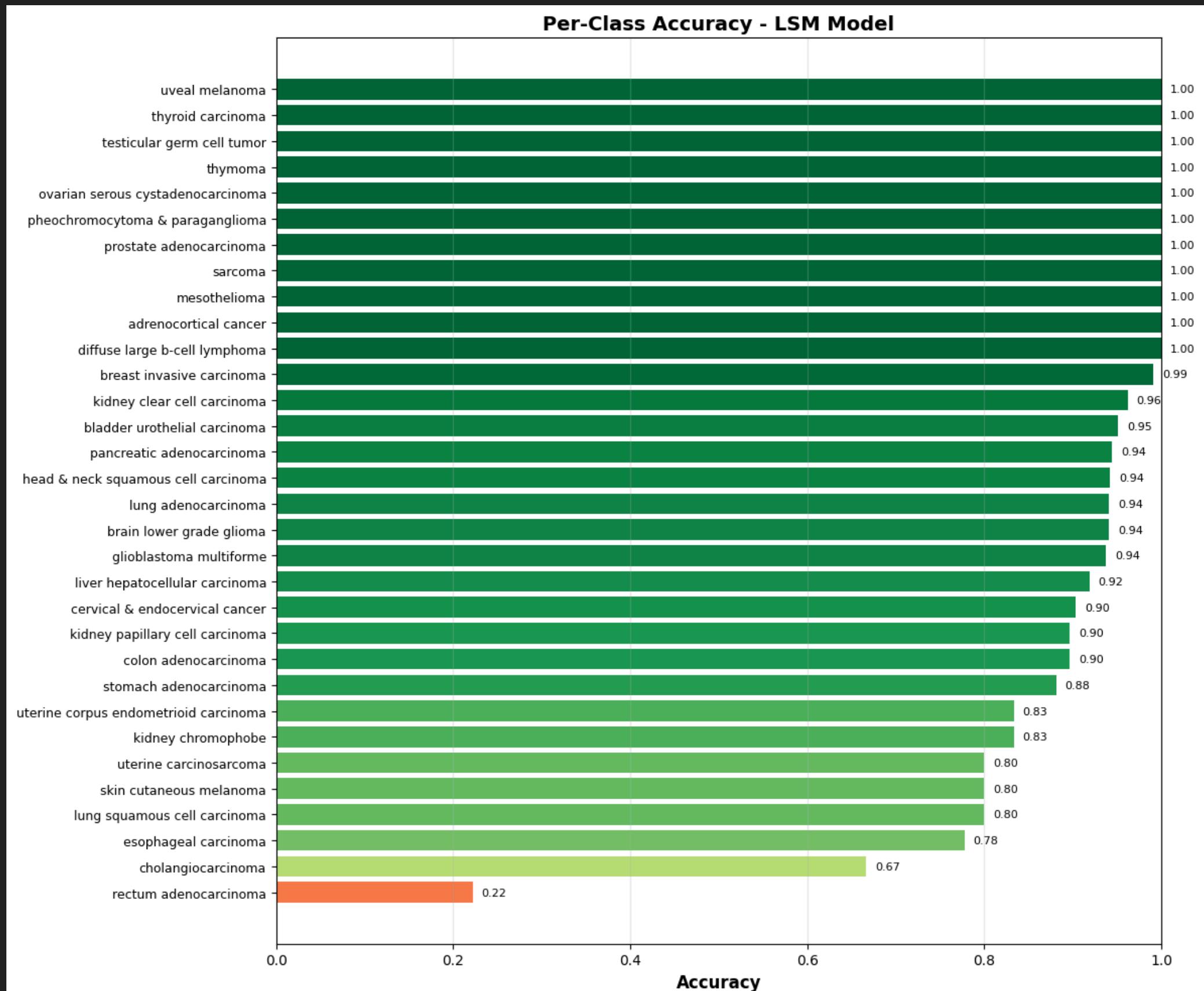


Comparison

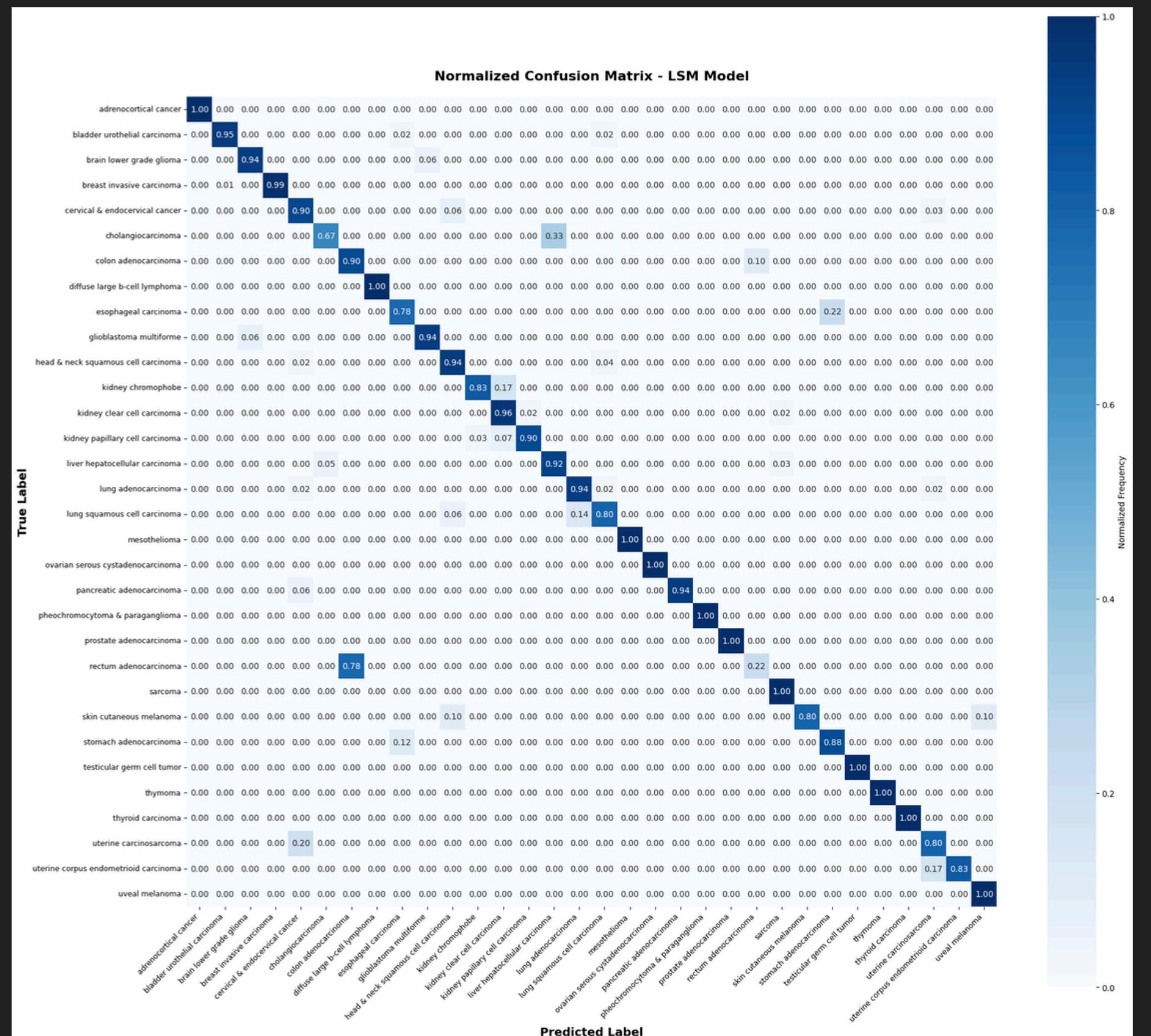
Results



Types



Confusion Matrix

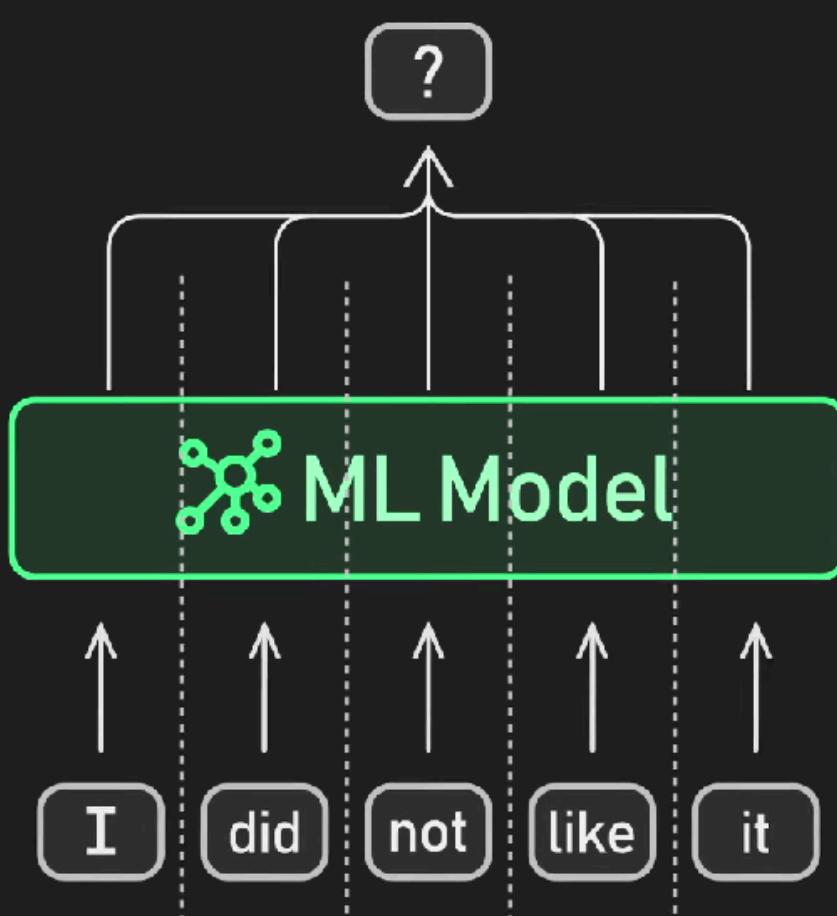


Transformers

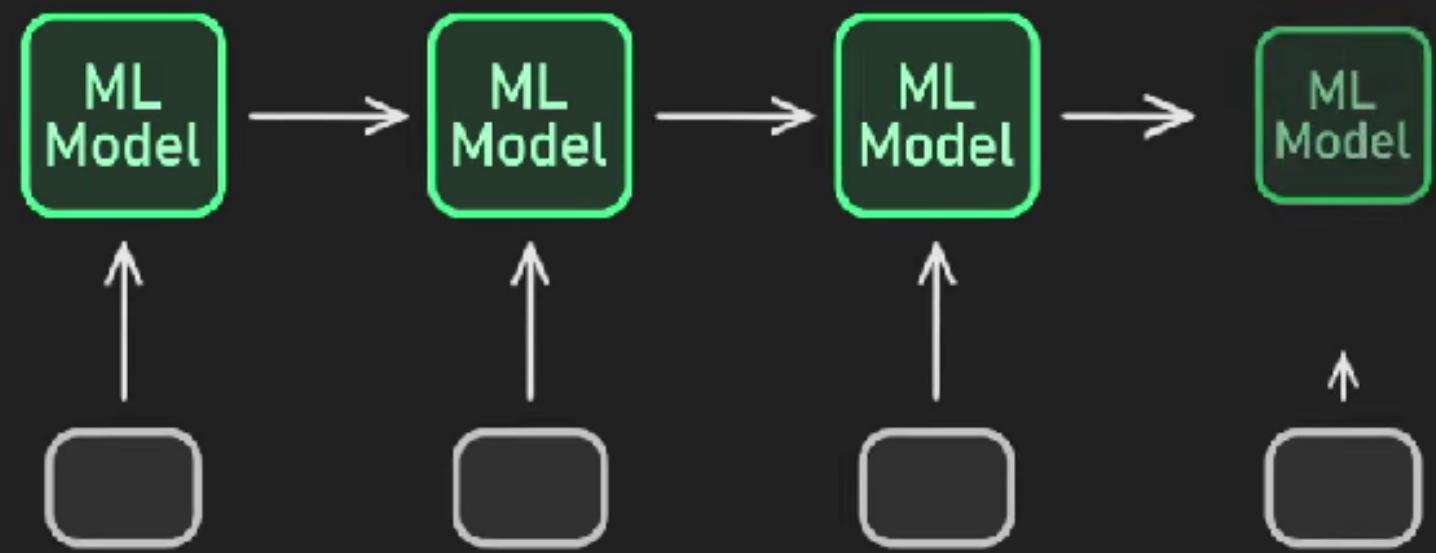


Before Transformer Problems

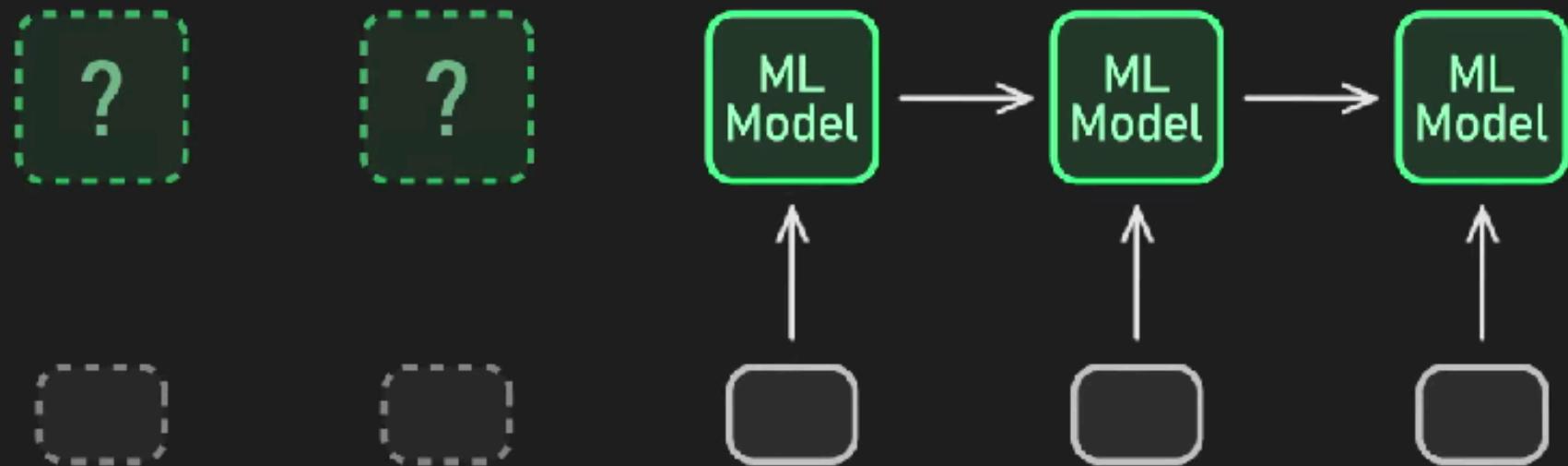
Tokens independent
Lost context



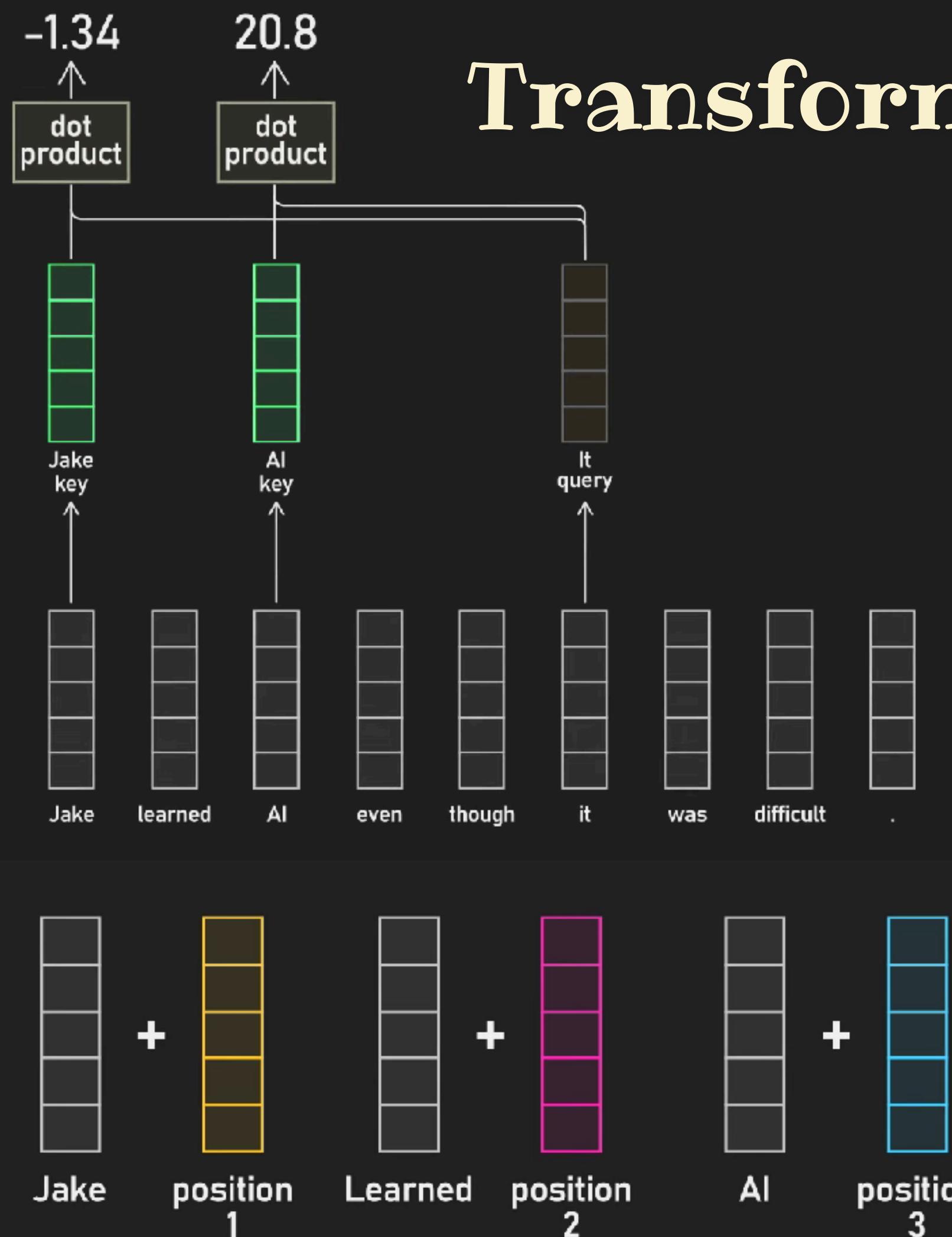
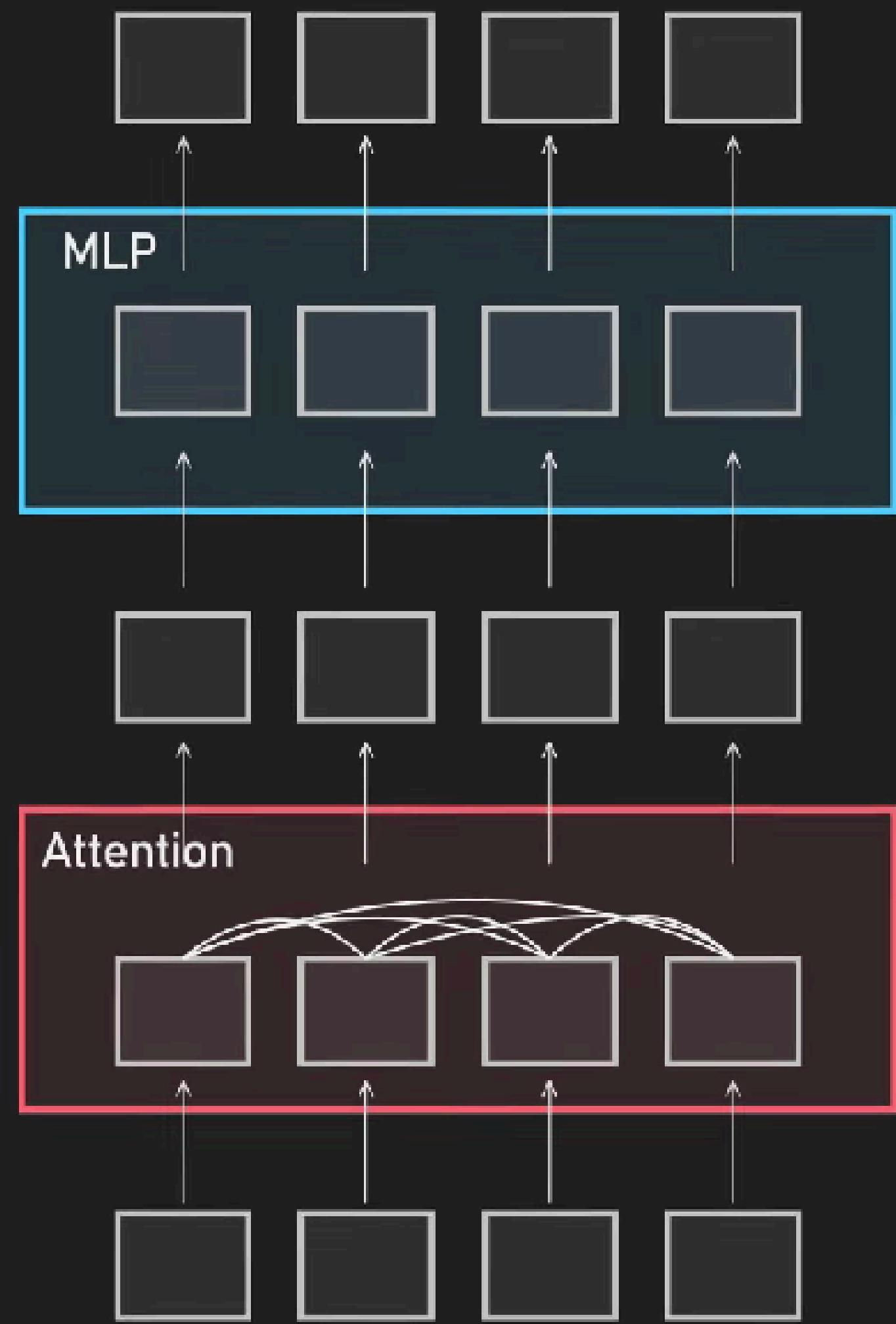
Sequential

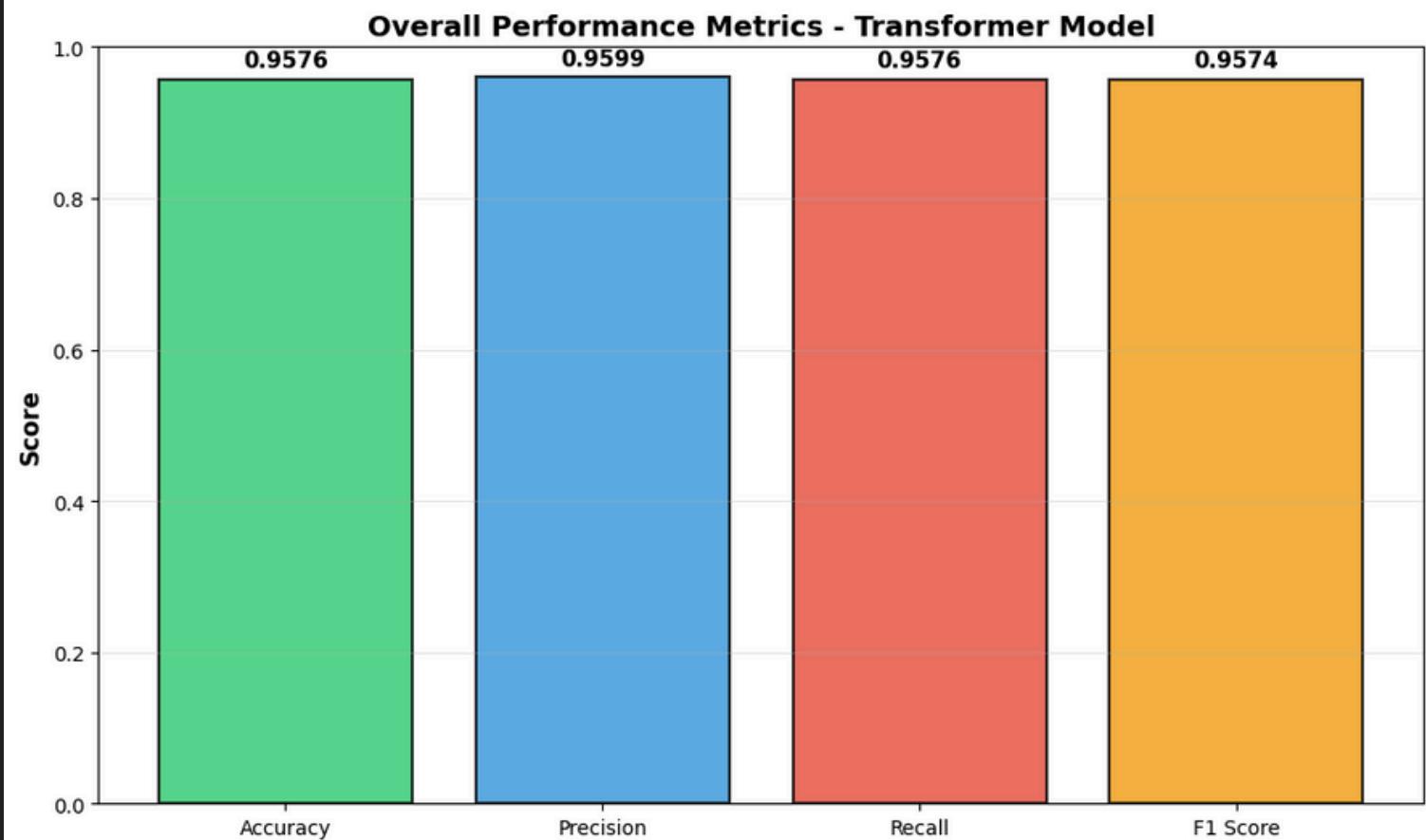


No Long Term
Dependencies



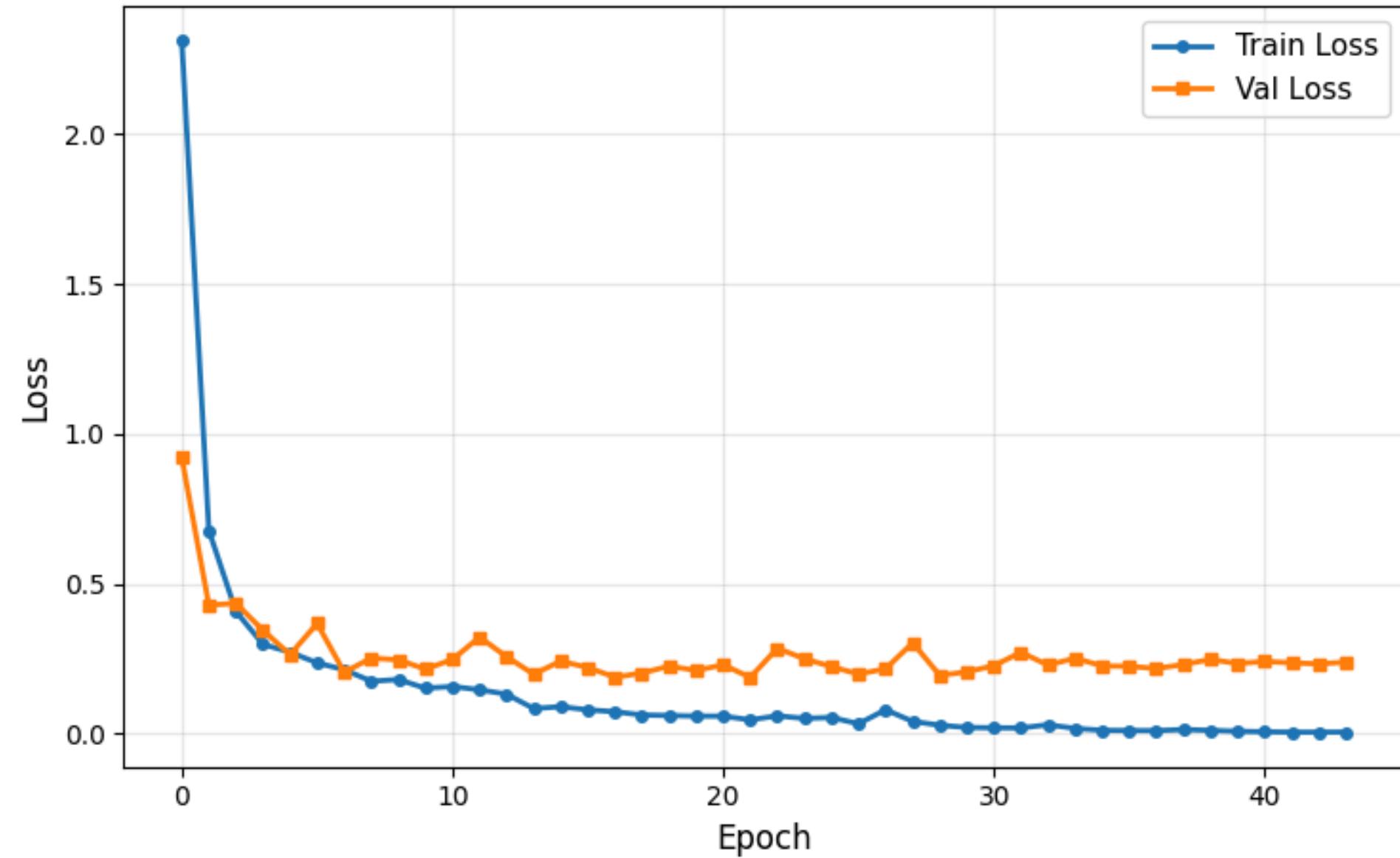
Transformer



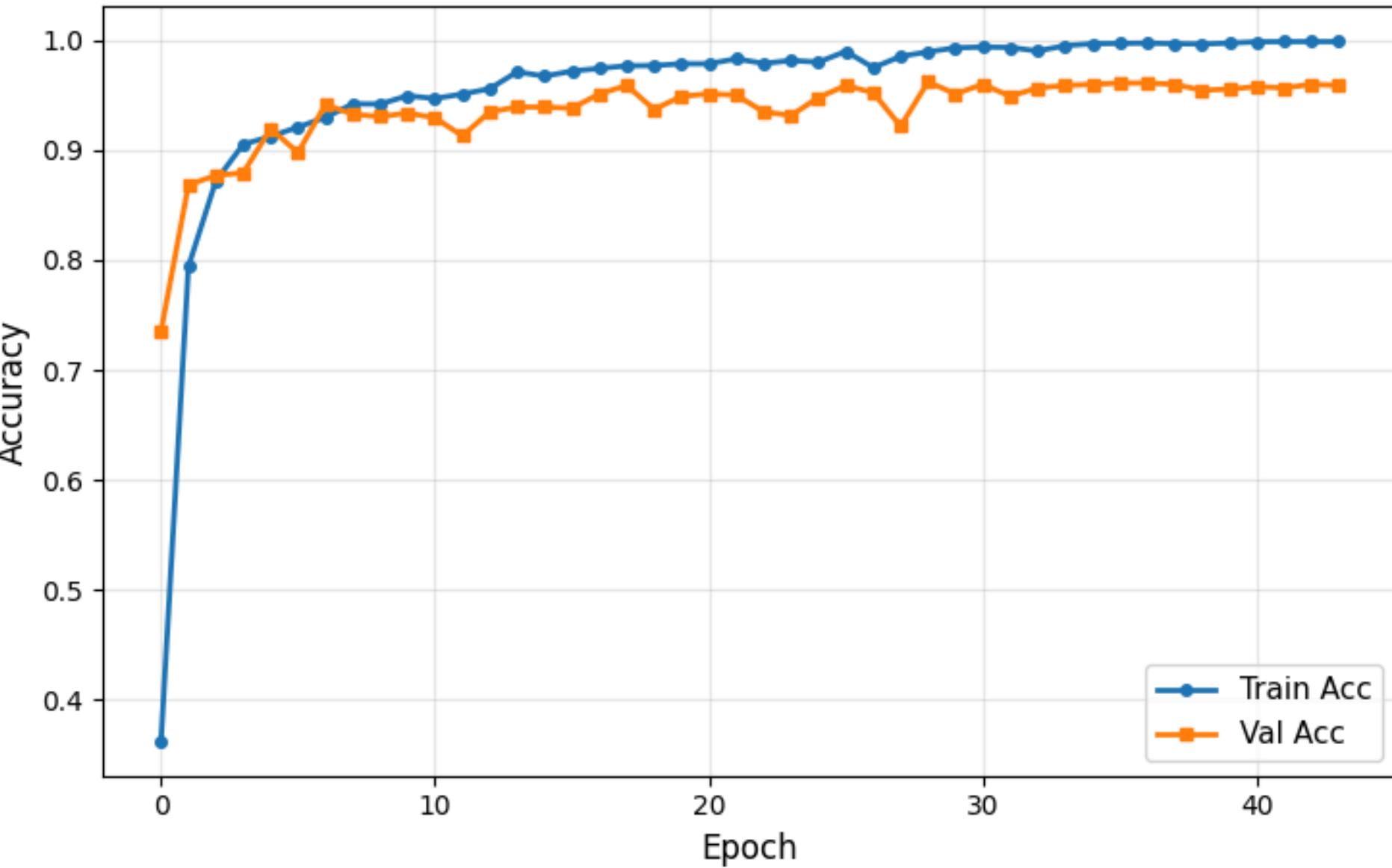


Training

Training and Validation Loss

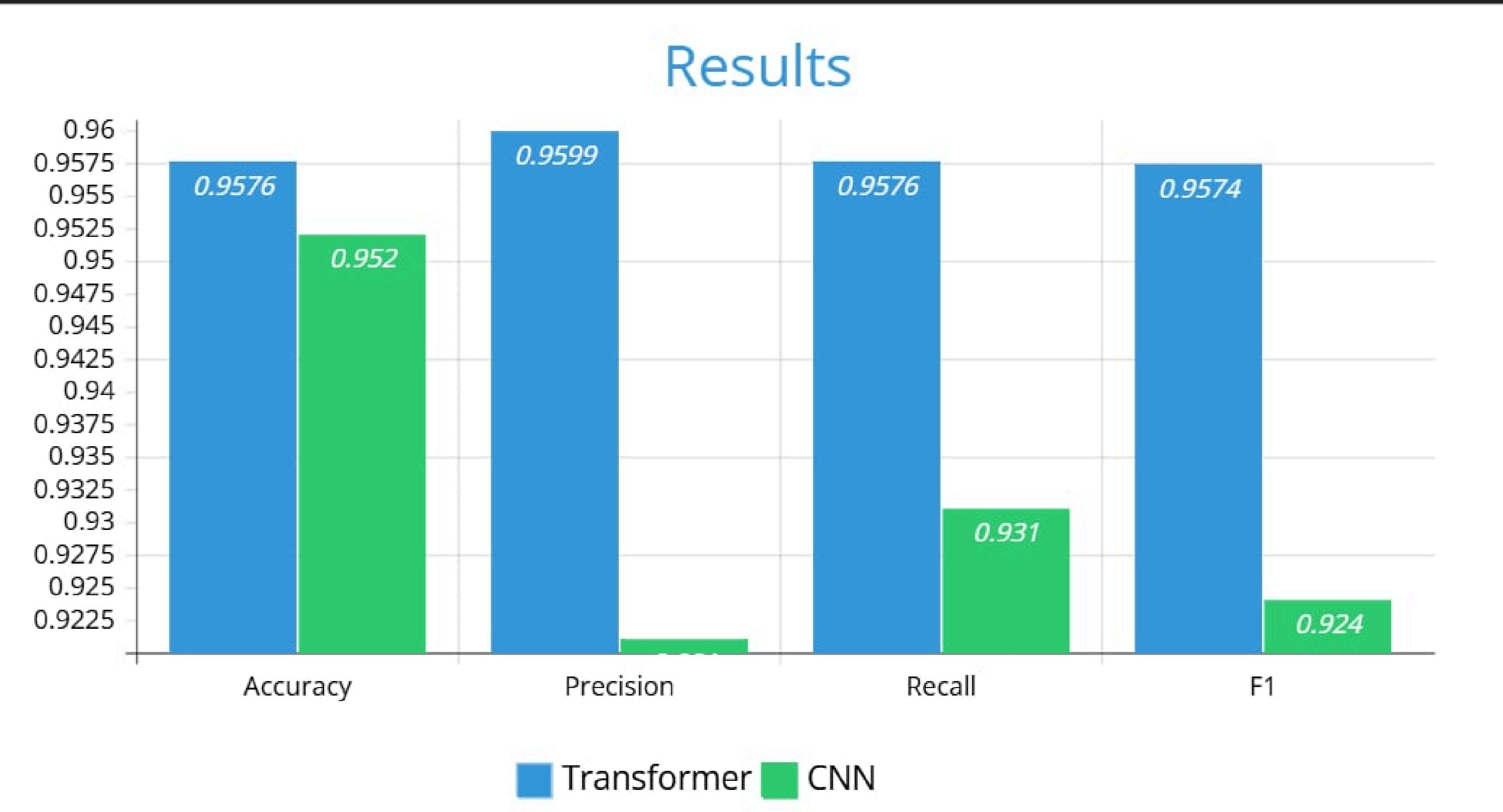


Training and Validation Accuracy

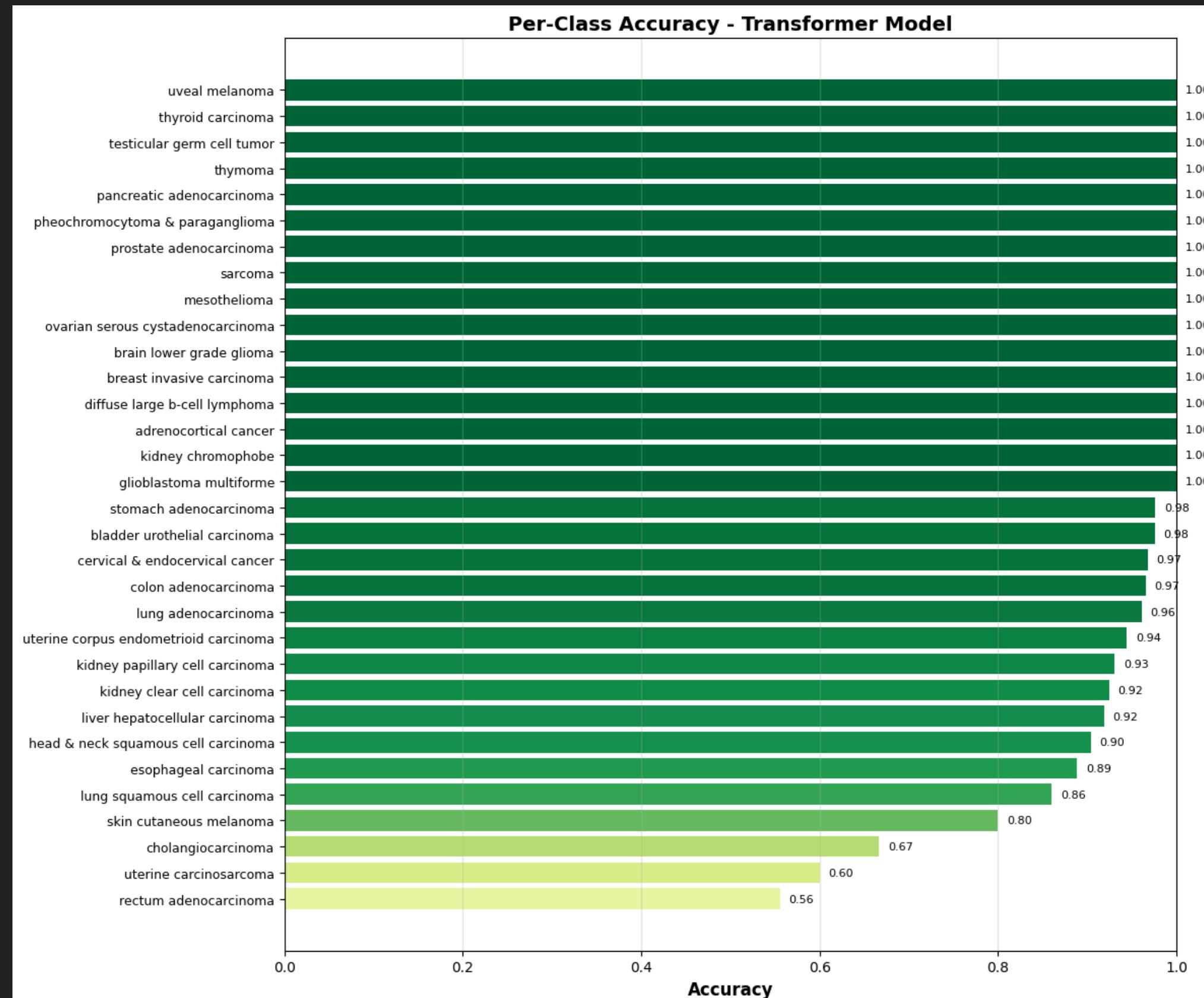


Comparison

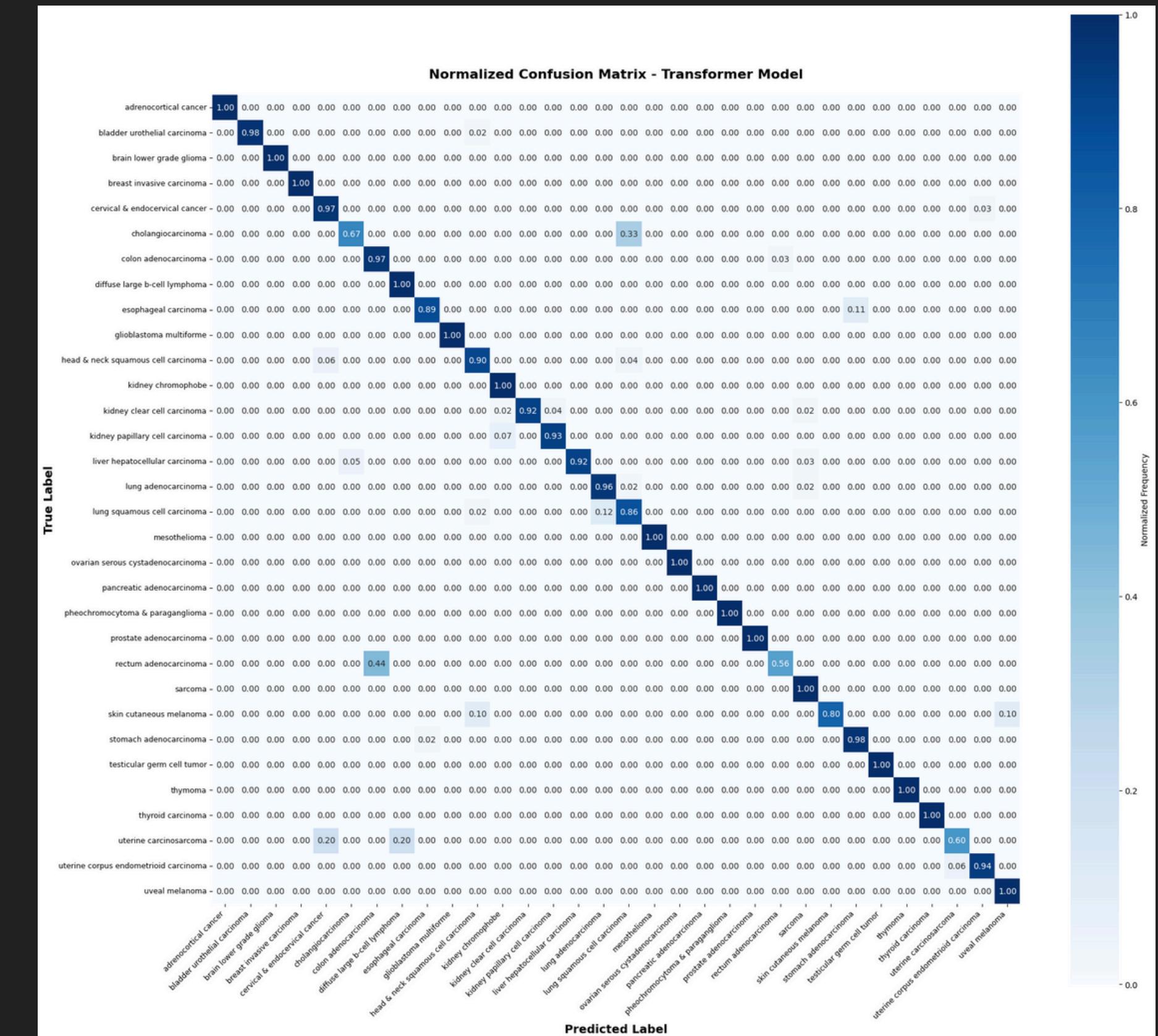
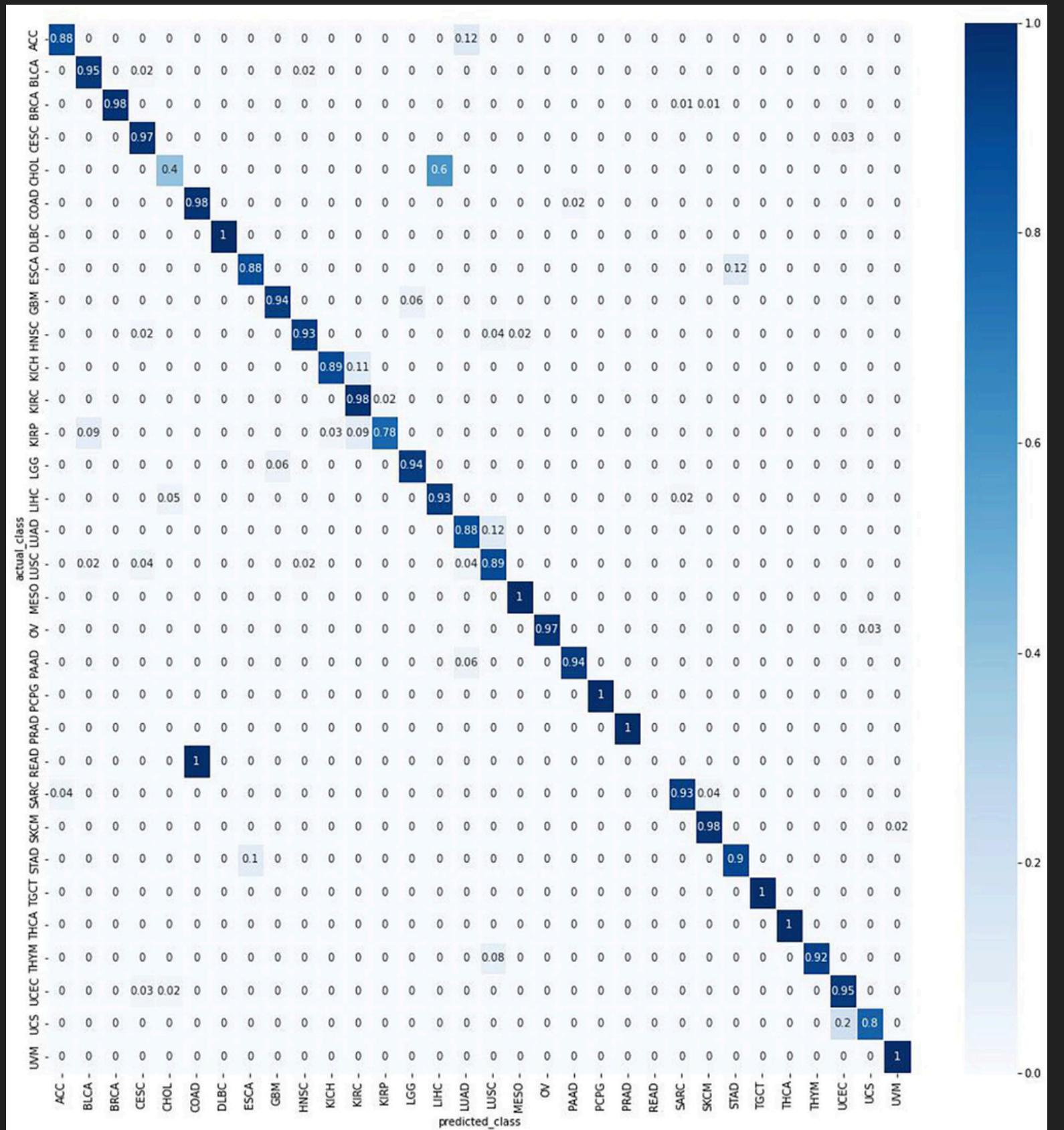
Results



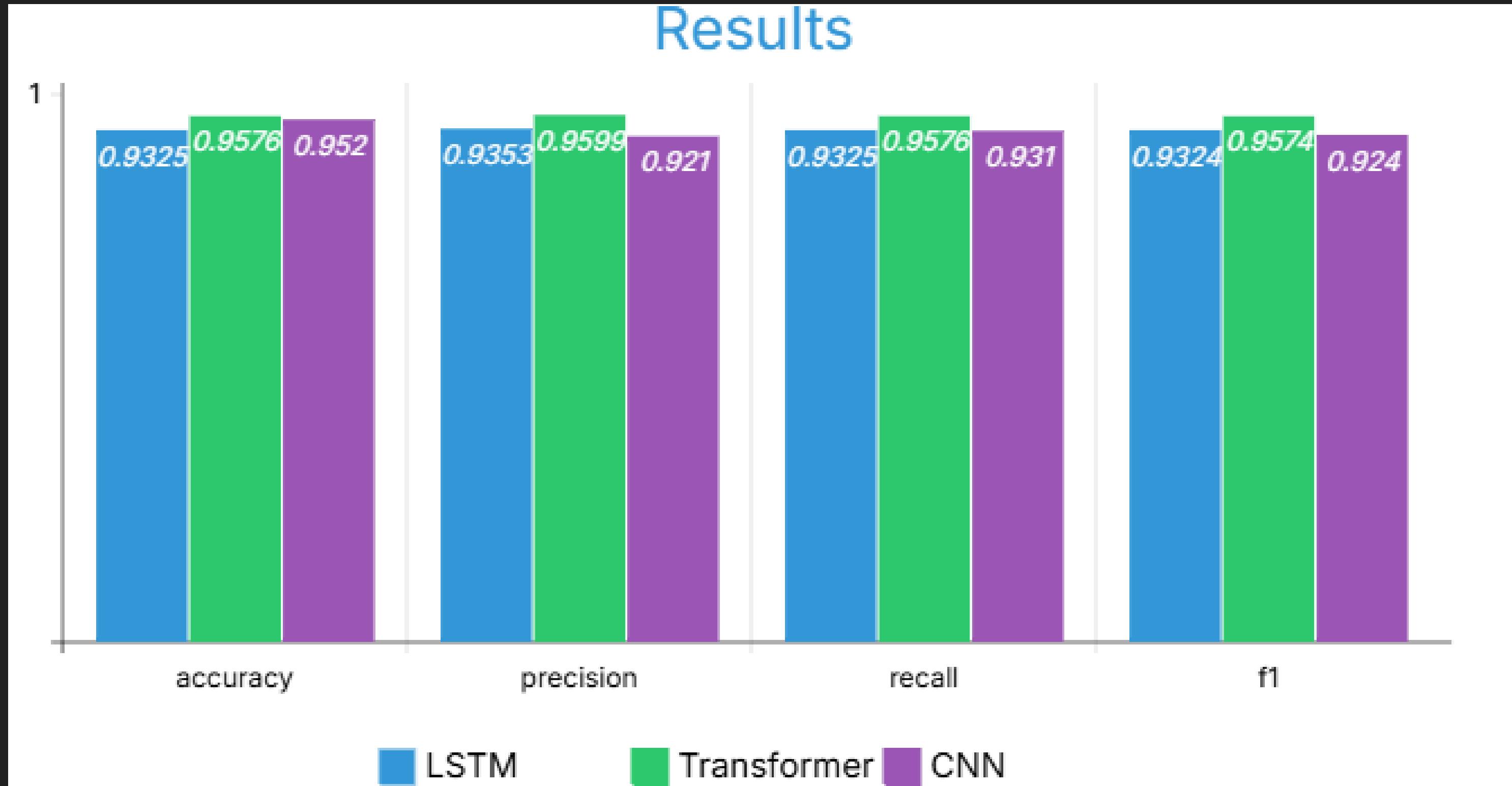
Types



Transformer VS CNN



Summary



Thank
You

