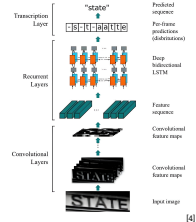


Úvod

OCR (Optical character recognition) je metoda, která rozpoznává tištěný či ručně psaný text a vrací jeho digitální zázpis. Pro tuto techniku se již celkem standardně využívají neuronové sítě. Starší architektury bývaly postaveny na kombinaci konvolučních a rekurentních neuronových sítí. Konvoluční neuronové sítě obvykle zastávají roli extraktoru znaků. Díky jejich správné utilitě je možné detekovat text na obrázcích s různým sořevěním a na odlišných pozicích. Rekurentní neuronové sítě se dále starají o překlad znaků na znaky. V kontextu neuronových sítí lze mluvit o architektuře encoder-decoder, kdy je encoder tvořen konvolučními vrstvami a decoder vrstvami rekurentními.

V poslední době jsou prováděny experimenty s novými moderními strukturami. Konvoluční neuronové sítě jsou nahrazovány vision transformery postavenými na attention mechanismu. Experimentuje se také s úplným vyloučením encoderu: výsledné neuronové sítě tak bývají postaveny na principu **decoder-only**. Jako dekodery se využívají například generativní předtrénované transformery (GPT), které vstupní obrázky kódují na tokeny a následně iterativně generují znaky až do rozpoznání konce testové sekvence. Generativní transformery se ukázaly jako výsoce efektivní napříč různými typy úloh, aktuálně jsou nejvíce používány v rámci velkých jazykových modelů (LLM).



Datové sady

V práci chceme dosáhnout co nejlepších výsledků na některém z ručně psaných datasetů. Žili jsme proto datasete LAM[1], který obsahuje zápisy Italského autora Ludovico Antonio Muratori, pořízené během 60 let jeho života. Dataset je rozdělen na stránky a řádky datasetu dle řádku. Dalším použitým datasetem je IAM, anglicky ručně psaný dataset rozdělený dle autora. Texty typicky obsahují malé množství slov (1–2).

Třetím využitým datasetem je námi vytvořený. Jedná se o útržky z italského korpusu (oscar), generované nástrojem Synthesizer, simulující ručně psaný text, s množstvím perturbací.

Datové sady jsou rozděleny v různých poměrech (okolo 80-10-10) na trénovací, validační a testovací sady.



TTOCR

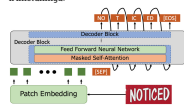
Model TtOCR[2] je postaven na architektuře encoder-decoder s využitím vizuálního transformeru DeiT pro zpracování obrázu a jazykového transforméru pro generování textu. Pro ručně psané texty jsme zvolili verzi small kvůli rychlosti a jednoduchosti. Při prvním finetuningu jsme dosáhli velmi špatných výsledků (CER kolem 0.9)



kvůli chybě dvojnásobku normalizaci vstupních dat. Po nápravě se sice přesnost zlepšila, ale model stále selhával – rozpoznával jen polovinu znaků, patrně kvůli škodlivým vstupním obrázkům na 384 × 384, což mohlo narušit čitelnost kurziv. Kvůli tímto omezením jsme se rozhodli pro jiný přístup.

DTtOCR

Architektura DTtOCR[3], navržená Maseem Fujitakem, je postavena na decoder-only přístupu s využitím modelu GPT-2 od OpenAI. Jelikož parametry původního modelu nejsou veřejně dostupné, bylo nutné ho v rámci práce natrénovat od začátku. Na rozdíl od počátečního tréninku TtOCR jsme k problému přistoupili systematicky: model nejprve rozbíráme předtrénovanými latinskými texty, než jsme přistoupili k finetuningu.



Implementace

V rámci implementace jsme se rozhodli využít existující řešení DTtOCR a trénovat s využitím CTC loss a SAM optimizací. Inspirace na zmíněné úpravy pochází z [4]. Váha našeho projektu je především v experimentech nad zmíněnou architekturou, která je poměrně nová (2023) a stojí za přehledným prozkoumáním. Nemalou zásluhu na implementaci je generativní syntetický datový soubor, který se odlišuje dle cílové úlohy (ručně psaný text, tištěný text, umělecký text...).

Model byl nejprve předtrénován na 1.2 milionech vstupů ze syntetického datového souboru problémů jedné úlohy. Pro generování předtrénování byla následně detekována nad LAM datasetem. Bylo provedeno 30 epoch nad trénovací datovou sadou, trénování bylo zastaveno, aby se předešlo overfittingu. Z výsledků implementace lze odvodit závěry o důležitosti variability a velikosti datové sady pro předtrénování a výkonu testované architektury.

CTC loss

(Connectisttemporal Classification)

Používání v seq-to-seq úlohách pro signální úlohy. V našem případě byla použita pro predikci sekvence znaků bez nutnosti explicitního zarovnávání mezi vstupní a výstupní sekvencí.

SAM

(Sharpness Aware Minimization)

Mýlenkové použití SAM na trénování modelu DTtOCR je snaha o zvýšení odolnosti vůči malým změnám, které mohou při detekci znaků nastat. Model by měl být schopen lépe generalizovat napříč různými podmínkami, druhy textů, etc.

Metriky

Nejdůležitějšími sledovanými metrikami při trénování modelu byly hodnota ztrátové funkce při validaci a hodnoty CER (Character Error Rate) a WER (Word Error Rate) zvláště během předtrénování modelu.

S počet substitucí (špatně rozpoznávaných znaků nebo slov)

D počet deletek (chybějících znaků nebo slov)

I počet insertů (nepřítušných znaků nebo slov)

N počet znaků v referenčním textu (pro CER)

W počet slov v referenčním textu (pro WER)

$$CER = \frac{S + D + I}{N}$$

$$WER = \frac{S + D + I}{N}$$

$$Accuracy = \frac{CER_{train}}{CER_{test}}$$

$$F1 = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall}$$

Zhodnocení

	fine-LAM	fine-IAM	fine-Synth	pre-LAM	pre-IAM	pre-Synth
Total CER	0.0781	0.5321	0.1654	0.2865	0.1775	0.0309
Total WER	0.1181	0.6509	0.3103	0.4224	0.2253	0.0678
Accuracy	0.9219	0.3901	0.8465	0.1655	0.7631	0.9369
Average CER	0.9262	0.1337	0.2830	0.0800	0.4624	0.7062
Average CER	0.0952	0.0919	0.1786	0.2479	0.1175	0.0222
Median CER	0.0000	0.0000	0.0460	0.2500	0.0000	0.0000
Std Dev CER	0.1120	0.5842	0.3142	0.1799	0.2582	0.0700
Average WER	0.1108	0.6595	0.3514	0.3638	0.2253	0.0995
Median WER	0.0000	0.0000	0.1667	0.3750	0.0000	0.0000
Std Dev WER	0.1736	0.7430	0.6972	0.2801	0.4305	0.1737

↓ 1. trénování výtěr respektive nižší hodnota jsou lepší

Předtrénování prokázalo vysoký výkon a značnou chybovost na obou hlavních metrickách CER a WER. Následně doladení na reálné datové sady LAM však vedlo ke zhoršení těchto metrik, přičemž celková CER vzrostla na 0.0781, celková WER na 0.1181 a přesnost klesla na 0.9219.

Vyhodnocení naznačuje, že krok doladení nezlepšil robustní základ vytvořený syntetickým předtrénováním napříč všemi datovými sadami. Nicméně lze tvrdit, že hlavní cíl byl splněn, neboť krok z předtrénování modelu na validaci je známý.

Z grafu nejběžnějších chyb při rozpoznávání slov lze obecně pozorovat tendenci, kdy model výslovně počítá samohlásky (např. „a“, „e“) nebo koncové hlásky v krátkých, běžně užívaných italských slovech.



Model má také tendenci dle více chyb čím dale postupuje v generování nebo rozpoznávání sekvence. Začátky sekvencí jsou rozpoznávány s větší přesností než jejich konec. Chyby se akumulují a model má potíže s udržením kontextu na delší vzdálenosti (error propagation).

Pro obě metriky (CER, WER) jsou zastoupeny nejvyšší frekvence u nízkých hodnot chybovosti, vzhledem k logaritmickeému měřítku osy y.



[1] Silvia Casanelli, Vittorio Pippi, Massimiliano Martin, Marcello Cornia, Lorenzo Baraldi, Kermorvan Christophe, and Rita Cucchiara. The lam dataset: novel benchmark for line-level handwritten text recognition. In International Conference on Pattern Recognition, 2022.

[2] Maseem Fujitake. Decoder-only transformer for optical character recognition, 2023.

[3] Minghao Li, Tenghao Lv, Lei Cui, Yipian Lu, Dinei Florencio, Cha Zhang, Zhongjun Li, and Furu Wei. Tracr: Transformer-based optical character recognition with pre-trained models, 2021.

[4] Rishu Rai, Sushiti Shinde, Pratiksha Sur, Swagnali Kulkarni, and Jayashree Jadhav. Automatic license plate recognition using yolo-v5 and tesseract-ocr. International Journal of Innovative Research in Computer and Communication Engineering, 10(10):656, 01 2008.

[5] Yuting Li, Dexiong Chen, Tinglong Tang, and Xi Shen. Htr-v: Handwritten text recognition with vision transformer. Pattern Recognition, 118:10967, 2025.