

Handwritten text recognition (OCR)

Bc. Pavel Raur (xraurp00), Bc. Jiří Žilka (xzilka11), Bc. Marián Zimmermann (xzimme03)*

Abstrakt

Projekt se zabývá rozpoznáváním ručně psaného písma pomocí OCR (Optical Character Recognition) s použitím self-supervised učení. V rámci projektu byly porovnány různé postupy self-supervised trénování OCR a následného doladování takto natrénovaného modelu. Byla zvolena architektura sítě pro následné experimentování s těmito metodami a byli provedeny experimenty. Následně tímto postupem bude vytvořeno několik kandidátních řešení pro porovnání.

Klíčová slova: OCR — self-supervised učení — strojové rozpoznávání ručně psaného textu

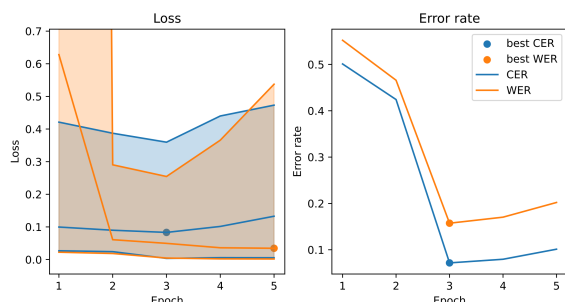
Přiložené materiály: [Stáhnutelný Kód](#)

*xraurp00@stud.fit.vut.cz, xzilka11@stud.fit.vut.cz, xzimme03@stud.fit.vut.cz,
Fakulta informačních technologií, Vysoké učení technické v Brně

1	1. Úvod	
2	Tato práce se zabývá natrénováním OCR (Optical	výřezu je jediný řádek textu. Dataset obsahuje celkem
3	Character Recognition) modelu pro dataset historick-	1153035 vzorků textu, z toho pouze 1% dat je anoto-
4	ého ručně psaného písma. Většinová část datasetu je	vaných (11469 vzorků). Anotovaná data jsou rozdělena
5	bez anotací. Projekt se zabývá zejména prozkoumáním	na trénovací (80%), testovací (7,5%) a validační sadu
6	metod self-supervised učení a prováděním experimentů	(12,5%).
7	k nalezení optimální kombinace metod, které povedou	
8	k co nejlepšímu výsledku na daném datasetu.	
9	2. Zvolený model a dataset	
10	Pro experimenty byl vybrán model <i>TrOCR</i> ¹ , který	2.2 Augmentace dat
11	používá transformerovou architekturu pro rozpoznávání	Pro zlepšení kvality modelu byla použita augmentace
12	textu na obrázcích a jeho přepis. Jak je popsáno v [3]	dat, což je běžný způsob využíváný při self-supervised
13	model je inovativní <i>state-of-the-art</i> řešení využívající	učení. Byly zvoleny následující augmentace:
14	transformery v oblasti, kde tradičně byly použity ze-	• Mixup
15	jména konvoluční neuronové sítě. Zároveň dosahuje	• Cutmix
16	kompetitivních výsledků s tradičními typy sítí. S ohle-	• Horizontální skosení
17	dem na rostoucí popularitu transformerů, ne jen v oblasti	• Natáhnutí/smrštění řádku
18	zpracování textu a obrazu, jsme se rozhodli právě pro	• Gaussian noise
19	tento typ architektury.	• Gaussian square - přidání šumu v podobě čtverce,
20	2.1 Dataset	který zakrývá náhodné místo o velikosti max
21	Použitý dataset je tvořený výřezů obrázků ručně psaného	40px
22	historického textu v anglickém jazyce. V každém	• Barevná maska
		3. Postup trénování
		Pro trénování byl použit model <i>microsoft/trocr-base-</i>
		<i>stage1</i> , který poskytuje základní váhy pro inicializaci
		TrOCR modelu. Váhy v tomto inicializačním modelu
		byly vybrány z modelů <i>BEiT</i> a <i>RoBERTa</i> pro inicial-
		izaci enkodéru pro detekci textu v obraze a dekodéru

¹<https://github.com/rsommerfeld/trocr>

47 pro přepis do textové podoby.[2] Tento model byl
 48 trénován 3 epochy na anotovaných datech. Ako vidno
 49 na 1 tak po 3 epochách dosahoval nejlepších výsledky.



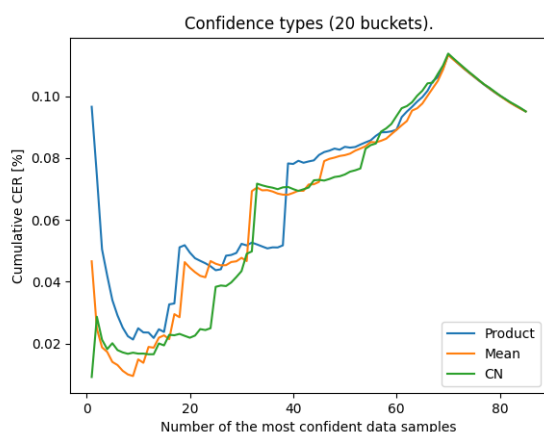
Obrázek 1. Vývoj CER a WER

50 3.1 Self-supervised učení

51 Tento model, který byl natrénován, byl následně použit
 52 k vytvoření anotací pro neanotovanou část datové sady.
 53 K tomu bylo nutné vybrat vhodnou metriku. Porov-
 54 nali jsme tři metriky a získali jsme pro ně následující
 55 plochy pod křivkou:

- 56 • Confidence product - 5.937035306340725
- 57 • Confidence mean - 5.643832876624217
- 58 • Confidence confusion network - 5.480604598787505

59 Na základě těchto výsledků jsme zvolili confusion net-
 60 work jako metriku. Na grafu 2 můžeme vidět, jak moc
 61 jsou chybové pseudo-labely při výběru dat pomocí této
 62 metriky.



Obrázek 2. Porovnání metrik

63 Pomocí zvolené metriky confusion network² se
 64 postupně vybíralo určité procento nejlepších anotací (10%,
 65 20% ... 100%) a přidali se do datasetu.

²<https://github.com/DCGM/pero-ocr>

Výchozí model	množství použitých labelů	počet epoch trénování
---------------	---------------------------	-----------------------

Tabulka 1. Caption

4. Vyhodnocování kvality modelu

Pro vyhodnocování modelů jsme použili CER (Char-
 acter Error Rate) a WER (Word Error Rate).[1] Vzorec
 pro výpočet:

$$CER = \frac{S + D + I}{S + D + I + C} \quad (1)$$

Kde:

- S(Substitutions) = počet správných znaků za-
měněných za chybné,
- D(Deletions) = počet chybějících znaků,
- I(Insertions) = počet znaků vložených navíc,
- C(Correct) = počet správných znaků.
- (toto je normalizovaná podoba CER)

Analogicky, WER je vypočítán nad slovy, na místo
 znaků. Někdy se udávají v podobě přesnosti (Ac-
 curacy), konkrétně CAR (Character Accuracy Rate)
 a WAR (Word Accuracy Rate).

$$CAR = 1 - CER \quad (2)$$

$$WAR = 1 - WER \quad (3)$$

5. Porovnání výsledků

Baseline model TrOCR Stage 1 byl před trénováním
 vyhodnocen na všech anotovaných datech a dosáhl
 následujících hodnot:

- CER = 0.3048
- WER = 0.5001

Model TrOCR Handwritten dosáhl hodnot:

- CER = 0.1768
- WER = 0.3524

Náš model, natrénovaný supervised příspupem z
 TrOCR stage1 dosáhl následujících výsledků:

- CER = 0.0717
- WER = 0.1573

Tento model byl následně použit pro generování pseudo-
 labelů pro další etapu trénování.

[[Doplnit data do tabulky!]] Modely natrénované
 na strojově anotovaných datech pak dosáhly následu-
 jících výsledků:

100	6. Kdo co udělal	Literatura	115
101	Pavel Raur	[1] LEUNG, K. Evaluate OCR Output Quality with Character Error Rate (CER) and Word Error Rate (WER). <i>Towards Data Science</i> . 2021. Dostupné z: https://towardsdatascience.com/evaluating-ocr-output-quality-with-character	116
102	• Výběr základního modelu pro trénování a použité architektury OCR.		117
103	• Základ kódů pro trénování modelu.		118
104	• Porovnání a selekce metrik pro výběr pseudolabelů.		119
105			120
106			121
107	Marián Zimmerman	[2] LI, M., LV, T., CUI, L., LU, Y., FLORENCIO, D. et al. <i>TrOCR: Transformer-based Optical Character Recognition with Pre-trained Models</i> . 2021. Dostupné z: https://huggingface.co/microsoft/trocr-base-stage1 .	122
108	• Výběr typů augmentací dat.		123
109	• Implementace augmentací.		124
110	• Dokumentace projektu.		125
111	Jiří Žilka	[3] LI, M., LV, T., CUI, L., LU, Y., FLORÊNCIO, D. A. F. et al. TrOCR: Transformer-based Optical Character Recognition with Pre-trained Models. <i>CoRR</i> . 2021, abs/2109.10282. Dostupné z: https://arxiv.org/abs/2109.10282 .	126
112	• Trénování OCR.		127
113	• Ukládání statistik.		128
114	• Dokumentace projektu.		129