

VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

Fakulta elektrotechniky
a komunikačních technologií

BAKALÁŘSKÁ PRÁCE

Brno, 2023

Martin Rosa



VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

BRNO UNIVERSITY OF TECHNOLOGY

FAKULTA ELEKTROTECHNIKY A KOMUNIKAČNÍCH TECHNOLOGIÍ

FACULTY OF ELECTRICAL ENGINEERING AND COMMUNICATION

ÚSTAV TELEKOMUNIKACÍ

DEPARTMENT OF TELECOMMUNICATIONS

SUPERROZLIŠENÍ V OBRAZE PRO ZAJIŠTĚNÍ VYLEPŠENÉHO MONITOROVÁNÍ ZABEZPEČENÝCH PROSTORŮ

SUPER RESOLUTION IN THE IMAGE TO ENSURE IMPROVED MONITORING OF SECURED AREAS

BAKALÁŘSKÁ PRÁCE

BACHELOR'S THESIS

AUTOR PRÁCE

AUTHOR

Martin Rosa

VEDOUCÍ PRÁCE

SUPERVISOR

doc. Ing. Radim Burget, Ph.D.

BRNO 2023

Bakalářská práce

bakalářský studijní program **Informační bezpečnost**

Ústav telekomunikací

Student: Martin Rosa

ID: 230648

Ročník: 3

Akademický rok: 2022/23

NÁZEV TÉMATU:

Superrozlišení v obraze pro zajištění vylepšeného monitorování zabezpečených prostorů

POKYNY PRO VYPRACOVÁNÍ:

Seznamte se s problematikou super-rozlišení, kde se zaměřte nikoli jen na využití ke zvýšení rozlišení fotografie, ale na využití v oblasti biometrie. Provedte podrobnou rešerši stavu vědy a techniky v této oblasti, zmapujte existující datové množiny pro tyto účely a vyberte vhodné metriky pro hodnocení kvality. S použitím vybraných architektur neuronových sítí natrénujte několik modelů, které poté s pomocí vybraných metrik vzájemně srovnajte v tabulce. Na vybraných příkladech ukažte případy, kdy metody selhávají a kdy naopak fungují.

DOPORUČENÁ LITERATURA:

podle pokynů vedoucího práce

Termín zadání: 6.2.2023

Termín odevzdání: 26.5.2023

Vedoucí práce: doc. Ing. Radim Burget, Ph.D.

doc. Ing. Jan Hajný, Ph.D.
předseda rady studijního programu

UPOZORNĚNÍ:

Autor bakalářské práce nesmí při vytváření bakalářské práce porušit autorská práva třetích osob, zejména nesmí zasahovat nedovoleným způsobem do cizích autorských práv osobnostních a musí si být plně vědom následků porušení ustanovení § 11 a následujících autorského zákona č. 121/2000 Sb., včetně možných trestněprávních důsledků vyplývajících z ustanovení části druhé, hlavy VI. díl 4 Trestního zákoníku č.40/2009 Sb.

ABSTRAKT

Cieľom bakalárskej práce bolo porovnať modely super-rozlíšenia s aplikáciou na reštaurovanie fotiek ľudských tvárí. V práci sme spracovali rešerš technológií superrozlíšenia a následne sme natrénovali a porovnali 5 modelov. Zameriavame sa hlavne na oblasť superrozlíšenia, ktorá by mohla byť nápomocná na identifikáciu osôb z bezpečnostných kamier. Použité technológie boli preto vyberané na základe percepčnej kvality a schopnosti identifikácie osoby na výstupnom snímku. Práca ukázala účinnosť porovnaných modelov pomocou objektívnych aj subjektívnych metrík. Výsledky boli porovnané v dotazníku (106 respondentov). Dotazník ukázal účinnosť použitia vlnovej transformácie v superrozlíšení tvárí.

KLÚČOVÉ SLOVÁ

superrozlíšenie, halucinácia tvárí, neurálne siete, hlboké učenie, bezpečnostné kamery, identifikácia

ABSTRACT

The point of this bachelor thesis was to compare models of super-resolution with the application on resolving human faces. A brief review of the technologies of super-resolution was created and five models were trained and compared. The focus was on the area of super-resolution that could be helpful with identifying people from CCTV cameras. Used technologies were therefore chosen based on their perceptual quality and ability to identify the person in the output image. This thesis has shown the effectivity of the compared models using objective and subjective metrics. The results were compared in a survey (106 respondents). Survey has shown the advantage of using wavelet-transform in the area of the super-resolution of human faces.

KEYWORDS

super-resolution, face hallucination, neural networks, deep learning, CCTV cameras, identification

ROSA, Martin. *Superrozlišení v obraze pro zajištění vylepšeného monitorování zabezpečených prostorů*. Brno: Vysoké učení technické v Brně, Fakulta elektrotechniky a komunikačních technologií, Ústav telekomunikací, 2023, 51 s. Bakalárska práca. Vedúci práce: doc. Ing. Radim Burget, Ph.D.

Vyhlásenie autora o pôvodnosti diela

Meno a priezvisko autora: Martin Rosa

VUT ID autora: 230648

Typ práce: Bakalárska práca

Akademický rok: 2022/23

Téma záverečnej práce: Superrozlišení v obraze pro zajištění vy-
lepšeného monitorování zabezpečených
prostorů

Vyhlasujem, že svoju záverečnú prácu som vypracoval samostatne pod vedením vedúcej/cého záverečnej práce, s využitím odbornej literatúry a ďalších informačných zdrojov, ktoré sú všetky citované v práci a uvedené v zozname literatúry na konci práce.

Ako autor uvedenej záverečnej práce ďalej vyhlasujem, že v súvislosti s vytvorením tejto záverečnej práce som neporušil autorské práva tretích osôb, najmä som nezasiahol nedovoleným spôsobom do cudzích autorských práv osobnostných a/alebo majetkových a som si plne vedomý následkov porušenia ustanovenia § 11 a nasledujúcich autorského zákona Českej republiky č. 121/2000 Sb., o práve autorskom, o právach súvisiacich s právom autorským a o zmene niektorých zákonov (autorský zákon), v znení neskorších predpisov, vrátane možných trestnoprávných dôsledkov vyplývajúcich z ustanovenia časti druhej, hlavy VI. diel 4 Trestného zákonníka Českej republiky č. 40/2009 Sb.

Brno

.....

podpis autora*

*Autor podpisuje iba v tlačenej verzii.

POĎAKOVANIE

Rád by som sa poďakoval vedúcemu bakalárskej práce pánovi doc. Ing. Radimu Burgetovi, Ph.D. za odborné vedenie, konzultácie, trpezlivosť a podnetné návrhy k práci.

Obsah

Úvod	11
1 Superrozlíšenie	12
1.1 Prevzorkovanie	12
1.1.1 Interpoláčn� metody	12
1.1.2 Prevzorkovanie zalo�en� na u�en�	13
1.2 Single-image vs. multi-image	14
2 Aplik�cia neur�lnych siet� na superrozl�šenie	16
2.1 Stratov� funkcie	16
2.2 Konvolu�n� neur�ln� siete	17
2.3 Superrozl�šenie zalo�en� na GAN sie�ach	18
2.4 Architekt�ra typu Transformer	20
2.4.1 Transformery v oblasti superrozl�šenia	22
3 Superrozl�šenie tv�r�	24
3.1 Unik�tne v�zvy FSR	24
3.1.1 Kvalitat�vne a kvantitat�vne meranie kvality	24
3.1.2 Probl�my GAN	24
3.2 Superrozl�šenia s vopred zn�mymi inform�ciami	25
3.2.1 Orienta�n� body	25
3.2.2 Tepeln� mapy	25
3.2.3 Parsovacie mapy	25
3.2.4 Atrib�ty tv�re	26
3.2.5 Inform�cia o identite	26
3.3 Superrozl�šenie so zachovan�m identity	26
3.3.1 Zachovanie identity pomocou rozpozn�vania tv�re	27
3.3.2 Zachovanie identity p�rov�mi metodami	27
4 Popis experimentu	28
4.1 D�tov� sady	28
4.1.1 CelebA	28
4.1.2 MLFDB	29
4.2 Metody merania kvality	29
4.2.1 PSNR	29
4.2.2 SSIM	29
4.2.3 LPIPS	30
4.3 Zvolene technol�gie	30

4.3.1	EIPNet	30
4.3.2	WaveletSRNet	33
4.3.3	NLSN	34
4.3.4	HAT	35
4.3.5	SRDD*	37
5	Výsledky experimentu	40
5.1	Dotazník	42
	Záver	44
	Literatúra	45
	Zoznam symbolov a skratiek	51

Zoznam obrázkov

1.1	Vrstva Sub-Pixel. Zdroj [40] (prekreslené)	14
2.1	Ukážka úspechu GAN (v zátvorkách sú hodnoty PSNR a SSIM). Zdroj [22] (upravené)	19
2.2	Schéma modelu transformer. Zdroj [44] (prekreslené)	21
2.3	Schéma modelu transformeru s využitím na superrozlíšenie. Zdroj [2] (prekreslené)	23
3.1	Ukážka vopred známych informácií. Zdroj [13] (preložené)	26
4.1	Nákres hranového bloku. Zdroj [17] (prekreslené)	31
4.2	Ukážka fotky z modelu EIPNet	32
4.3	Ukážka fotky z modelu WaveletSRNet	34
4.4	Ukážka fotky z modelu NLSN	35
4.5	Ukážka rozdielu zachytenia obrysov obrázku po aplikácii sobelovho operátora pri rozdielnych veľkostiach	36
4.6	Ukážka fotky z modelu HAT	37
4.7	Porovnanie aktivácií na sieti SRCNN s pomocou ReLU a s pomocou xUnit. Zdroj [18] (preložené)	38
4.8	Schéma pozornostnej brány. Zdroj [35] (prekreslené)	38
4.9	Ukážka výstupu modelu SRDD*	39
5.1	Kvalitatívne porovnanie výsledkov. V zátvorkách sú PSNR [dB], SSIM a LPIPS nasledovne.	41
5.2	Výsledky dotazníku	42

Úvod

Pod pojmom superrozlíšenie sa rozumie, zväčšovanie rozlíšenia obrázku, tak aby došlo k zvýšeniu jeho kvality. Super-rozlíšenie je téma zo širokou škálou aplikácie (bezpečnosť, zdravotníctvo, atď.). V našej práci sa budeme zaoberať použitím týchto technológií na bezpečnosť. Konkrétne riešime aplikáciu super-rozlíšenia na vylepšovanie záberov z bezpečnostných CCTV kamier pre účely identifikácie osôb na záberoch z týchto kamier.

Tento problém je komplexný a existuje veľa metód, ako sa ho pokúsiť vyriešiť. Jeho riešenie je však podstatné pre jeho významnú aplikáciu pre bezpečnostné zložky, ktoré pravidelne pracujú s zábermi na základe ktorých sú identifikovaný podozrivý z trestných činov. Zvýšená kvalita týchto záznamov (ktoré zvyknú často trpieť veľmi zlou kvalitou) môže identifikáciu uľahčiť, príp. ju umožniť. V tejto práci porovnávame moderné technológie super-rozlíšenia pri ich aplikácii na rekonštrukciu obrazu ľudských tvárí. Výsledky hodnotíme pomocou objektívnych metrík a subjektívnych metrík, za účelom zistenia aplikovateľnosti týchto technológií na daný problém.

V úvode práce je stručný popis problému super rozlíšenia a jeho aplikácia na tváre. Následne sú popísané technológie zvyšovania kvality obrazu. V ďalších kapitolách sú popísané neurónové siete a ich aplikácia na danú problematiku. Zameriavame sa tu aj na unikátne problémy, ktoré sa objavujú v prípade použitia superrozlíšenia za účelom identifikácie. Na záver uvádzame technológie použité pri experimentácii.

Hlavným prínosom práce je porovnanie 5 technológií superrozlíšenia objektívnymi aj subjektívnymi metrikami. Zároveň sme upravili existujúcu sieť za účelom dosiahnutia lepších výsledkov. Nami upravenej sieti sa podarilo dosiahnuť výsledky porovnateľné s modernými technológiami, v subjektívnom hodnotení skončila na druhom mieste. V subjektívnom hodnotení sme ukázali účinnosť vlnovej transformácie pri superrozlíšení ľudských tvárí.

Dielo končí časťou kde uvádzame výsledky práce. Zaoberáme sa tu spôsobmi ich hodnotenia a procesom tréningu sietí. V závere zhrnieme subjektívne aj objektívne výsledky.

1 Superrozlíšenie

Pod pojmom superrozlíšenie (super-resolution) sa rozumie zvýšenie kvality snímky, prípadne videa. Superrozlíšenie má veľa oblastí (superrozlíšenie tvárí, satelitných snímok, medicínskych snímok a pod.) v tejto práci sa zameriavame na superrozlíšenie tvárí. Kľúčovou časťou superrozlíšenia je tzv. prevzorkovanie (upsampling).

1.1 Prevzorkovanie

Superrozlíšenie je založené na zmene snímky s nízkym rozlíšením (LR) na snímku z vysokým rozlíšením (HR). V skutočnosti ide o od základu zle založený problém, keďže z viacerých HR snímok je možné spraviť rovnakú LR snímku. Toto znemožňuje získanie pôvodnej HR snímky, výsledkom superrozlíšenia je len aproximácia HR snímky.[46] Pri získavaní HR snímky musíme zvýšiť rozlíšenie pôvodnej LR snímky. Matematicky možno problém superrozlíšenia zapísať ako:

$$\widehat{HR} = F(LR, \theta), \quad (1.1)$$

kde LR reprezentuje LR snímku, F reprezentuje zvolenú metódu superrozlíšenia a θ reprezentuje parametre funkcie. \widehat{HR} reprezentuje aproximáciu pôvodného HR obrázku. Parametre funkcie závisia od použitej technológie. Pri superrozlíšení založenom na neurálnych sieťach sa získavajú iteratívnym optimalizačným procesom nazývaným učenie. [46]

Metódy superrozlíšenia sú založené na prevzorkovaní. Jedná sa o pridávanie pixelov do obrázku na základe daných pravidiel. V nasledujúcich odsekoch sa zameriame na rôznorodé metódy prevzorkovania.

1.1.1 Interpolačné metódy

Jedná sa o najznámejšie a najjednoduchšie metódy prevzorkovania. Patrí sem interpolácia najbližších susedov (nearest neighbor interpolation), bilinéarna interpolácia (bilinear interpolation), bikubická interpolácia (bicubic interpolation) a iné.

Interpolácia najbližších susedov

Jednoduchý algoritmus, založený na princípe najbližších susedov. Pridané pixely si hodnoty preberú od pôvodných pixelov, ktoré sú k nim najbližšie. Jedná sa o veľmi rýchlu a výpočetne nenáročnú metódu. Metóda však produkuje výsledky s veľmi nízkou kvalitou. [46]

Bilineárna interpolácia (BIL)

Štatistická metóda. Doplnené pixely si hodnoty získavajú vypočítaním “priemerných” hodnôt pôvodných pixelov. Jednoduchá, rýchla a nenáročná metóda. Dosahuje lepšej kvality, ako interpolácia najbližších susedov. [46]

Bikubická interpolácia (BIC)

Pracuje na rovnakom princípe ako BIL. Pridáva komplexnosť, keďže, narozdiel od BIL, vypočítava hodnoty pridaným pixelom pomocou polynomickej metódy. “Priemerné” hodnoty sa počítajú na krivke, čo dosahuje lepšie výsledky, ako BIL. Jedná sa o najpopulárnejšiu metódu interpolácie. [46]

1.1.2 Prevzorkovanie založený na učení

Vyššie popísané metódy používajú iba informácie získané zo snímky na zvýšenie rozlíšenia. Momentálne však dosahuje lepších výsledkov prevzorkovanie založené na učení, ktoré používa aj informácie naučené použitou neurálnou sieťou.[6]

Konvolučné neurálne siete

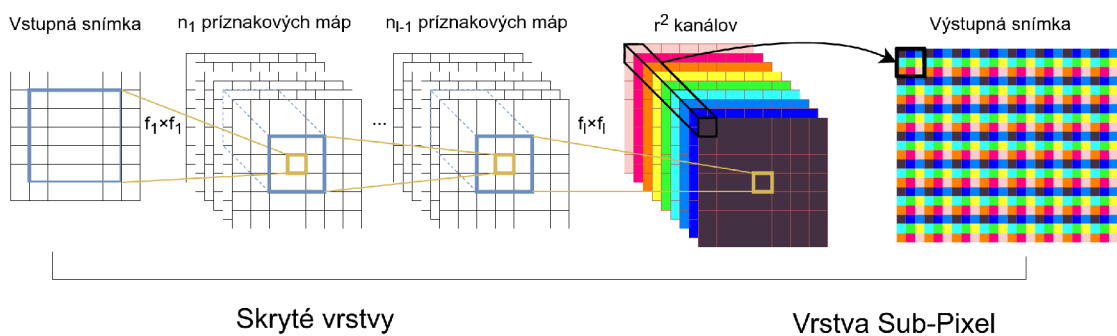
Konvolučné neurálne siete priniesli veľký pokrok v oblasti superrozlíšenia.[10] [6] Táto technológia je založená na tzv. konvolúciách. Konvolúcie sú proces, kde po snímke prechádzame kernelom s danou veľkosťou a posunom (stride). Kernel (niekedy nesprávne nazývaný aj filter) sa pohybuje po snímke, pričom rýchlosť posunu je určená posunom, a sčítava hodnoty, ktoré zachytáva. Výsledkom je tzv. Skonvolovaná vlastnosť (convoluted feature). Jedná sa o snímku so zníženým rozlíšením, ktorá sa následne posiela do ďalších vrstiev siete. Tento proces umožňuje sieti zachytiť podstatné vlastnosti snímky. Čo sa využíva hlavne v problémoch počítačového videnia (computer vision).[10]

Transponovaná konvolúcia

Tento mechanizmus funguje na opačnom princípe, preto je často nesprávne nazývaný dekonvolúcia. Pri transponovanej konvolúcii najprv snímku rozšírime na požadovanú veľkosť, pričom pridané pixely doplníme nulovými hodnotami. Následne do snímky opäť doplníme pixely s nulovými hodnotami, tak aby sme dosiahli dvojnásobok požadovanej veľkosti. Na záver, po snímke necháme prejsť kernel, ktorý snímku skonvoluje na požadovanú veľkosť. Týmto spôsobom dokáže sieť na základe príznakových máp (feature maps) predpovedať, ako by mohla *HR* snímka vyzeráť. Táto metóda sa, vďaka svojim uspokojivým výsledkom používa v mnohých technológiách. [46] [10]

Vrstva Sub-pixel (Sub-pixel layer)

Vrstva Sub-pixel vytvára viaceré kanály pomocou viacnásobnej konvolúcie. Konkrétne sa pôvodný obrázok konvoluje, až pokiaľ nevznikne r^2 kanálov, kde r je faktor zväčšenia. Pri vstupe s rozmermi $W \times H \times C$ (kde W je šírka snímky, H je výška snímky a C je počet farebných kanálov) bude teda výstupom $W \times H \times C * r^2$. Následne je na výstup aplikovaná tzv. zmiešavacia (shuffle) operácia, kde dôjde k prevodu $W \times H \times C * r^2$ na $rW \times rH \times C$. Týmto spôsobom dôjde ku zvýšeniu rozlíšenia o faktor r . Túto metódu môžeme vidieť bližšie na 1.1. Veľkou výhodou tejto operácie je, že pri konvoluci a následnom upsamplingu nie je nutné pridávať nuly, ani obrázok nejakým spôsobom umelo rozširovať. Vďaka tomuto si dokáže sieť získať viac kontextuálnych informácií, čo umožňuje lepšie generovať details. Zároveň, podobne ako pri Transponovanej konvolúcii, je operácia pomerne výpočtetne jednoduchá. Prvým príkladom siete, ktorá túto metódu používala je ESPCN, ktorá dosiahla lepšie výsledky ako vtedy dovtedy existujúce technológie a zároveň to dokázala v lepšom výpočtovom čase. [40]



Obr. 1.1: Vrstva Sub-Pixel. Zdroj [40] (prekreslené)

1.2 Single-image vs. multi-image

Superrozlíšenie môžeme rozdeliť na dve kategórie

- Superrozlíšenie z jedného snímku (Single-image super-resolution - **SISR**)
- Superrozlíšenie z viacerých snímok (Multi-image super-resolution - **MISR**)

Ako z názvu vyplýva, superrozlíšenie z jedného snímku vychádza z použitia jedného snímku ako základ, ktorému je následne zväčšené rozlíšenie. Zatiaľ čo pri superrozlíšení z viacerých snímok vychádzame z viacerých snímok, ktoré použijeme ako základ, pričom si z nich berieme informácie, podľa ktorých vygenerujeme *HR* snímku.

Nešťastnou realitou pre MISR je, že väčšina literatúry sa venuje práve SISR. Nejedná sa o dobre preskúmanú oblasť. MISR sa dostáva pozornosti iba v superrozlíšení satelitných snímok [15] [38]. V poslednej dobe sa však do pozornosti dostáva aj MISR tváří [37].

Rok 2019 priniesol technológiu DeepSUM. Model sa skladá zo siete na extrakciu črt z snímok. Na základe týchto črt potom ďalšia sieť vytvára tzv. registračné filtre, na základe ktorých sa registrujú príznakové mapy. V poslednej časti siete dochádza k splývaniu príznakových máp a generovaniu \widehat{HR} snímky. [33]

Technológia RAMS priniesla využitie reziduálneho príznakového bloku. Tento blok umožňuje modelu zamerať sa na vysokofrekvenčné komponenty snímky. Model dosiahol podobných výsledkov ako DeepSUM. [39]

Jedna z mála technológií superrozlíšenia tvarí bola prestavená v [37]. Autori svoju sieť založili na U-Net sieti (podobne ako ostatné technológie MISR). Následne bol pridaný tzv. GEU blok (gated enhanced unit). Narozdiel od väčšiny technológií, bol tento model ohodnotený nie len pomocou objektívnych metrik ale aj na základe dotazníku, v ktorom respondenti porovnávali výstupy rôznych sietí s pôvodnou HR snímokou a následne hlasovali za najpodobnejšiu. V danom dotazníku sa ukázala ako najlepšia práve táto technológia. Z čoho môže vyplývať, že sieť dokáže dobre rekonštruovať identifikačné črty.

StatNet je model, ktorý priniesol tzv. štatistický blok. Tento blok dokáže pretvoriť LR snímky na 6 príznakových máp, ktoré sa používajú v ďalších častiach modelu. Výhodou štatistického bloku je, že neobsahuje trénovateľné časti, čiže je výpočtovo efektívny. Zároveň je možné do tohto bloku dať hocikolko vstupných LR snímok. Čo je unikátna vlastnosť v oblasti MISR. [34]

Jednou z najnovších MISR technológií je QA-Net z roku 2022. Tento model používa tzv. mapy kvality (quality maps). Ich úlohou je odfiltrovať časti snímok, ktorých informačná hodnota bola znížená napr. šumom. Model vďaka tomu dokáže diskriminovať voči týmto častiam a získať väčšinu komponentov z kvalitnejších častí snímok. [23]

Je nutné zmieniť, že QA-Net bol vytvorený pre účely superrozlíšenia satelitných snímok. V tomto prípade dáva zmysel zamerať sa len na kvalitné časti snímok na získanie lepších informácií, z ktorých je možné generovať lepšie \widehat{HR} snímky. To sa aj empiricky potvrdilo. [23] Je pre nás náročné si predstaviť využitie tohto fenoménu na superrozlíšenie tvarí, a teda aj jeho aplikovateľnosť na túto prácu. Keďže mapy kvality by mohli diskriminovať voči častiam obsahujúcim dôležité identifikačné črty.

2 Aplikácia neurálnych sietí na superrozlíšenie

Neurálne siete sa ukázali ako veľmi účinná metóda riešenie problematiky superrozlíšenia. Dnes sa v podstate stretávame už len z ich aplikáciou na tento problém. V nasledujúcich kapitolách sa budeme venovať prístupom k superrozlíšeniu za pomoci neurálnych sietí.

2.1 Stratové funkcie

Pod pojmom stratová funkcia (loss function), niekedy tiež prekladaná ako chybová funkcia, rozumieme funkciu, ktorá si za vstup berie \widehat{HR} a HR a jej výstupom je odchýlka \widehat{HR} od HR . Jedná sa o nutnú súčasť všetkých technológií super-rozlíšenia založených na neurálnych sieťach. Stratová funkcia vedie neurálnu sieť na základe odchýlky predošlých výsledkov k lepším výsledkom, pričom úlohou siete je minimalizovať túto odchýlku. V nasledujúcej časti sa pozrieme na niektoré stratové funkcie. Bude tu použitá rovnaká notácia ako v 1.1. [46]

Funkcie pixlovej straty (Pixel-wise loss)

Funkcie pixlovej straty vypočítajú odchýlku na každom pixely. Skladajú sa z \mathcal{L}_1 ktorá vypočíta priemernú absolútnu chybu (mean absolute error) a \mathcal{L}_2 , ktorá vypočíta priemernú kvadratickú chybu (mean squared error)

$$\mathcal{L}_1(\widehat{HR}, HR) = \frac{1}{hwc} \sum_{i,j,k} |\widehat{HR}_{i,j,k} - HR_{i,j,k}|, \quad (2.1)$$

$$\mathcal{L}_2(\widehat{HR}, HR) = \frac{1}{hwc} \sum_{i,j,k} (\widehat{HR}_{i,j,k} - HR_{i,j,k})^2, \quad (2.2)$$

kde h, w, c sú výška, šírka a počet kanálov vstupných snímok nasledovne. [46]

Nevýhodou \mathcal{L}_2 je, že dobre nezachytáva malé rozdiely medzi snímkami. Výsledky sietí používajúce túto funkciu zvyknú byť preto príliš hladké bez zachovania menších detailov. V tomto má \mathcal{L}_1 značnú výhodu oproti \mathcal{L}_2 , keďže malé aj veľké rozdiely zachytáva podobne. [13] Ďalšou výhodou \mathcal{L}_1 je, že je inverzne korelované k jednej z najpoužívanejších hodnotiacich metrík – PSNR (peak signal-to-noise ratio). To znamená, že znižovaním \mathcal{L}_1 sa bude zároveň zvyšovať PSNR. [46]

Asi najväčšou nevýhodou funkcií Pixel-wise loss je, že nehodnotia oblasti, ako je napr. percepčná kvalita. Tento hendikep spôsobuje, že výsledky neurálnych sietí sú vzorovo púliš monotónne, bez náznaku po detailoch. [46]

Percepčná strata

Funkcie percepčnej straty fungujú na princípe zisťovania, na koľko boli zachované isté črty v \widehat{HR} oproti HR . Získavanie týchto črt sa robí pomocou na to prispôbenej neurálnej siete (napr. VGG [41]). Funkciu percepčnej straty možno zapísať, ako

$$\mathcal{L}_{\text{percept}}(\widehat{HR}, HR, \psi, l) = \|\psi^l(HR) - \psi^l(\widehat{HR})\|_2, \quad (2.3)$$

kde ψ značí zvolenú neurálnu sieť na extrahovanie črt. A l značí l -tú vrstvu siete.

Keďže percepčná stratová funkcia núti sieť vytvárať snímky s vyššou percepčnou kvalitou, snímky vyzerajú esteticky krajšie aj keď zvyknú mať nižšie PSNR oproti metódam využívajúcim Pixel-wise loss. [13]

2.2 Konvolučné neurálne siete

Konvolučné neurálne siete (convolutional neural networks – CNN) majú veľa foriem, no v základe sa skladajú z troch hlavných častí – konvolučná vrstva (convolutional layer), združovacia vrstva (pooling layer) a plne prepojená vrstva (fully-connected layer).

Konvolučná vrstva používa kernely na skonvolovanie vstupnej snímky na príznakové mapy. Veľa konvolučných sietí obsahuje viacero konvolučných vrstiev, tzn. tento proces prebieha viackrát. Každý neurón konvolučnej vrstvy je prepojený so špecifickými neurónmi predošlej konvolučnej vrstvy.

Združovacia vrstva slúži na zníženie rozlíšenia príznakových máp. Jej úlohou je dosiahnuť shift-invariance. Združovacie vrstvy sa väčšinou dávajú medzi jednotlivé konvolučné vrstvy.

Úlohou plne prepojenej vrstvy je prepojiť všetky neuróny predošlej vrstvy s momentálnou vrstvou. Táto vrstva nebýva vždy prítomná a môže sa nahradiť 1×1 konvolučnou vrstvou.

Poslednou vrstvou v CNN je výstupná vrstva (output layer). Štruktúra tejto vrstvy závisí na úlohe konkrétnej konvolučnej neurálnej siete. Napr. pri klasifikačných úlohách sa používa funkcia softmax. [10]

Konvolučné neurálne siete predstavujú jeden z prvých prelomov v oblasti super-rozlíšenia. Jednou z týchto sietí je SRCNN. Jedná sa o veľmi jednoduchú technológiu. Obsahuje iba 3 vrstvy, no napriek tomu dokáže generovať snímky s vyššou kvalitou, ako dovtedy používané metódy. [6]

V roku 2016 bola vymyslená DRCN. Táto technológia priniesla veľmi hlbokú rekurzívnu vrstvu, čo umožnilo vysoké zvýšenie kvality oproti doposiaľ existujúcim

technológiám. [16] Dong et al. sa pokúsili vylepšiť rýchlosť a znížiť výpočetnú obťažnosť SRCNN. Ich práca, FSRCNN, vymenila bikúbickú interpoláciu za transponovanú konvolúciu a pridala iný tvar konvolučnej siete. [7] O rok na bola vymyslená ESPCN. Tu bola prvý krát použitá vrstva Sub-pixel, čo umožnilo dosiahnuť lepšie kvantitatívne výsledky. [40] Ďalšou technológiou bola LapSRN. Tu bola vymyslená sieť, ktorá dokáže extrahovať črty z LR snímok, na základe ktorých sú vytvárané tzv. reziduálne snímky. Tieto snímky sa kombinujú s upsamplend snímkami na dosiahnutie väčšej kvality. Zároveň LapSRN dosahuje podobného výpočetného času a zložitosti ako FSRCNN. [20] Ying et al. vytvorili DRRN, ktorá používa rekurzívne učenie. Táto technológia dosiahla omnoho lepšie výsledky ako predošlé. [43] Následne vznikla EDSR kde boli použité reziduálne bloky, resp. ResNet čo umožnilo prehĺbenie siete s tým, že je možné sa vyhnúť degradácii výsledku. Táto technológia vyhrala NTIRE2017 výzvu. [25] Yiqun et al. Použili metódu nelokálnej riedkej pozornosti (Non-local sparse attention) na vytvorenie NLSN, ktorá dosahuje porovnateľných výsledkov. [30]

2.3 Superrozlíšenie založené na GAN sieťach

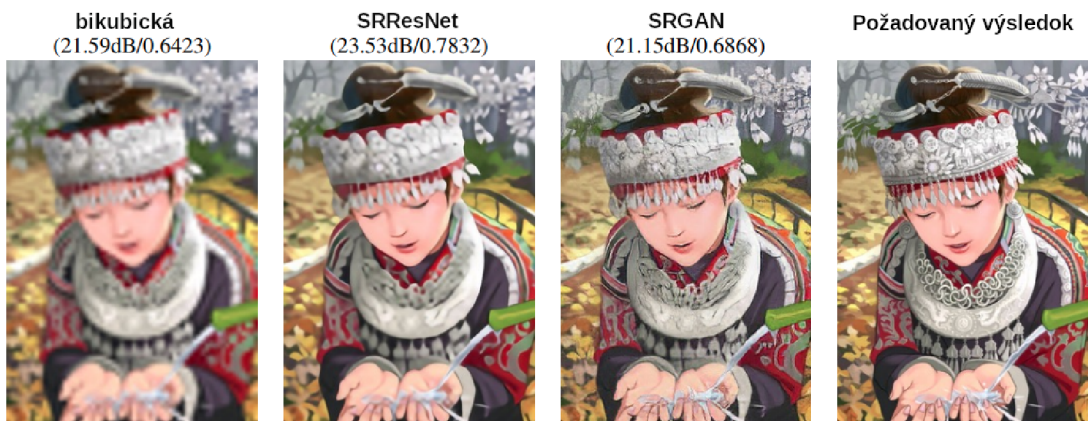
Rok 2014 priniesol nový pokrok v oblasti strojového učenia. Bol tým práve koncept generatívnej aversarialnej siete (generative adversarial network – GAN). Táto technológia zlepšila schopnosť umelej inteligencie generovať „svoje vlastné“ výtvary. GAN sú založené na princípoch genetických algoritmov. Konkrétne GAN sa skladá z dvoch neurálnych sietí.

Prvá sa nazýva generátor. Je to prevažne dekonvolučná sieť, ktorej úlohou je generovať chcené dáta na základe naučených dát. Generátor je nutné naučiť hodnotám reálnych dát, na základe ktorých chceme vygenerovať „falošné“ dáta, ktoré budú mať čo najbližšie k reálnym.

Diskriminátor, ktorý je druhou neurálnou sieťou v GAN, má za úlohu zistiť, či je výstup z generátora „dostatočne skutočné“. Diskriminátor je prevažne konvolučná sieť, ktorá je naučená na vstupoch reálnych dát, a na základe nich zisťuje s akou pravdepodobnosťou je výstup z generátora skutočný alebo umelo vytvorený. Inak povedané, pointou generátoru je oklamať diskriminátor tým, že vygeneruje dáta, ktoré budú tak podobné so skutočnými, že diskriminátor nebude schopný vidieť rozdiel medzi nimi.

Pôvodní autori tohto konceptu prirovnali GAN ku zápasu medzi falšovateľmi peňazí a políciou, ktorá sa ich snaží odhaliť (falšovatelia sú generátor a polícia je diskriminátor). Zároveň, ako postavy v spomínanej analógii sa aj generátor a diskriminátor dokážu zlepšovať vzájomným zápasením. [9]

V oblasti superrozlišenia zaznamenali GAN veľký úspech. Jednu z prvých implementácií je sieť SRGAN (super-resolution generative adversarial network). Ako generátor tu bola použitá sieť SRResNet. Jedná sa konvolučnú sieť, ktorá dokáže zvyšovať rozlíšenie snímok, zatiaľ čo diskriminátor sa snaží zistiť či je výsledok generátoru skutočná snímka alebo len prevzrokovaná. Vďaka tomuto SRGAN vytvára snímky z vysokou vizuálnou kvalitou. [22] O dva roky bola vytvorená ESRGAN. Táto technológia pridala RRDB bloky a upravila diskriminátor tak aby predpovedal skutočnosť snímok pomocou relatívnej skutočnosti, namiesto absolútnej. Pomocou týchto vylepšení, ESRGAN dokázal prekonať SRGAN v oblasti vizuálnej kvality, čím vyhral PIRM2018-SR výzvu. [45] V tom istom roku bola vytvorená URDGN. Ide sa o jednu z najznámejších GAN v oblasti superrozlišenia tváří. WGAN-GP použili tzv. Wassertein GAN aby spravili sieť stabilnejšou. FCGAN používa architektúru diskriminátoru aj generátoru založenú na U-Net sieťach. TDAE je technológia zameraná na superrozlišenie tváří s šumom a bez zarovňania. FaceAttr používa atribút tváre (viď 3.2.4) na zlepšenie kvality. Zatiaľčo modernejší Super-FAN používa tepelných máp (viď 3.2.2) ako vopred známu informáciu, ktoré si dokáže aj sám generovať. FSRFCH taktiež používa tepelné mapy, aby bola schopná presnejšie predvídať štruktúru tváre. [26] MLGE používa obrisy tváre vďaka čomu dosahuje výborné výsledky. [19] EIPNet používa obrisy tváre aj zachovanie identity. [17] SiGAN prináša iný spôsob zachovania identity (viď 3.3.2) vďaka čomu dosahuje výbornú percepčnú kvalitu. [11]



Obr. 2.1: Ukážka úspechu GAN (v zátvorkách sú hodnoty PSNR a SSIM). Zdroj [22] (upravené)

Adversariálna strata

Pre správne fungovanie generatívnych adversariálnych sietí je treba špeciálnych, tzv. Adversariálnych stratových funkcií. Adversariálne stratové funkcie sa skladajú z dvoch funkcií – stratová funkcia generátora (\mathcal{L}_G) a stratová funkcia diskriminátora (\mathcal{L}_D). Pôvodné technológie využívali stratové funkcie založené na krížovej entropii,

$$\mathcal{L}_G(\widehat{HR}) = -\log(\mathcal{D}(\widehat{HR})) \quad (2.4)$$

$$\mathcal{L}_D(\widehat{HR}, HR) = -\log(\mathcal{D}(HR)) - \log(\mathcal{D}(\widehat{HR})), \quad (2.5)$$

kde \mathcal{D} reprezentuje diskriminátor. HR reprezentuje požadovanú snímku a \widehat{HR} reprezentuje dosiahnutú snímku. [13]

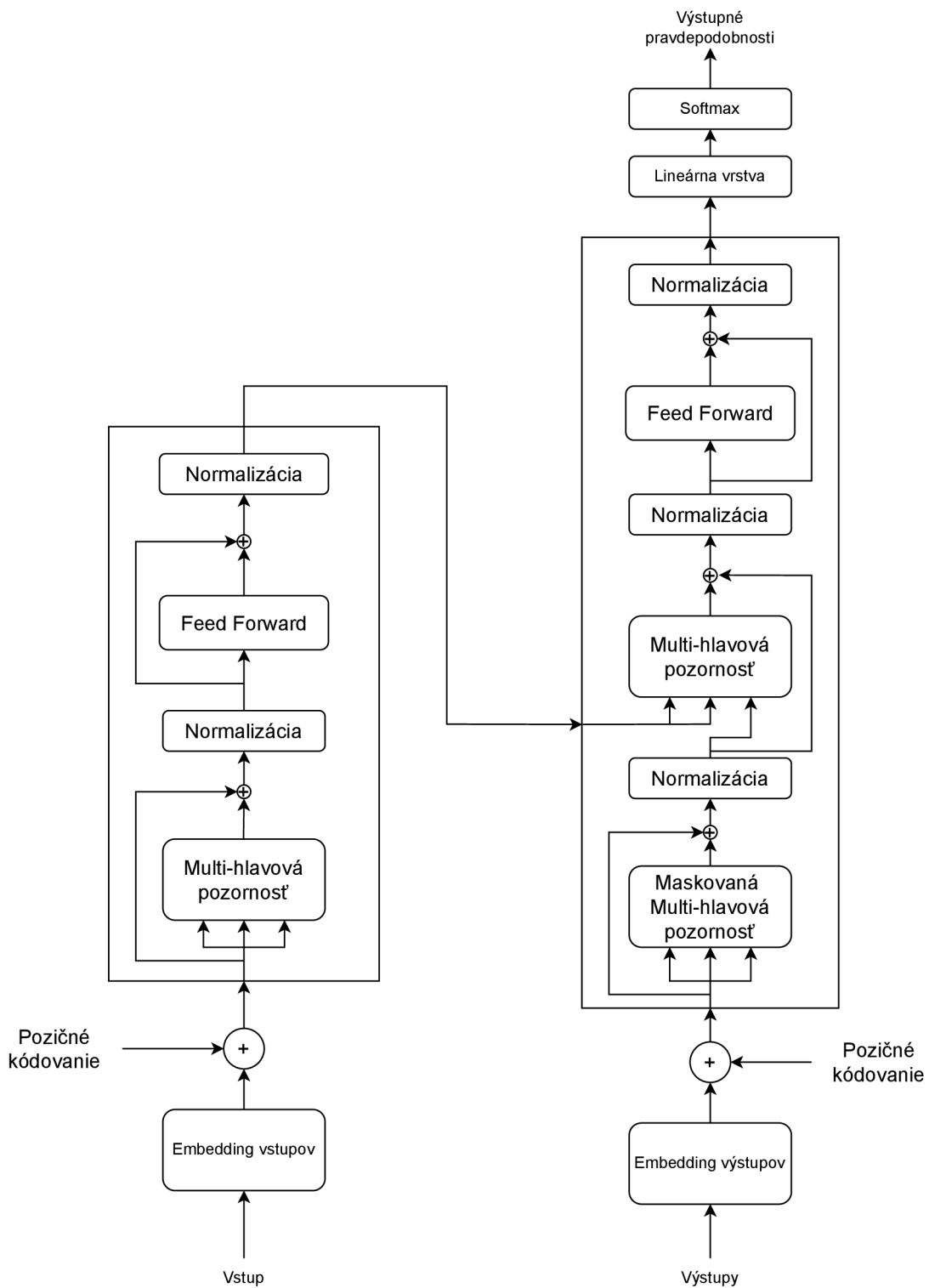
2.4 Architektúra typu Transformer

V roku 2017 bol vytvorený nový druh neurálnych sietí, zameraných na spracovanie jazyka. Transformer boli vytvorené za účelom náhrady komplikovaných rekurentných modelov za jednoduchšiu technológiu. Transformer sú založené na typickej štruktúre autoenkóderu. Ich netypickou častou je, že takmer výhradne používajú mechanizmy pozornosti. [44]

Transformer priniesol 3 hlavné príspevky:

- Pozičné kódovanie
- Pozornosť
- Seba-pozornosť

Pozičné kódovanie je spôsob akým dokážeme zaznamenať informáciu o pozícii slova vo vete. Pozornosťou sa rozumie mechanizmus akým sa model dokáže špeciálne zamerať na sémanticky podstatné dáta (v prípade spracovania jazyka to môžu byť kľúčové slová vo vete, v prípade počítačového videnia podstatné časti snímky). Mechanizmy pozornosti sa používali v neurálnych sieťach už aj pred vynálezom transformeru (hoci samozrejme nie v takej miere ako v transformeri). Unikátom transformeru je tzv. Seba-pozornosť.



Obr. 2.2: Schéma modelu transformer. Zdroj [44] (prekreslené)

Seba-pozornosť

Seba-pozornosť je nový druh pozornosťného mechanizmu. Výstupom seba-pozornosti je vážený priemer všetkých vstupov. Seba-pozornosť sa vždy aplikuje na dáta vychádzajúce z jedného predošlého bloku (preto sa nazýva seba-pozornosť). V transformeri je používaná tzv. multi-hlavovná pozornosť (multi-head attention). Multi-hlavová pozornosť je založená na škálovanej skalárne súčinovej pozornosti (scaled dot-product attention). Multi-hlavová pozornosť sa skladá z viacerých blokov, ktoré spolu bežia paralelne.

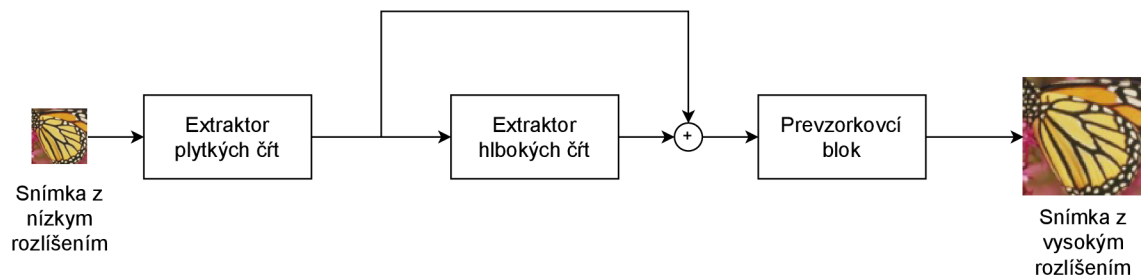
Veľkou výhodou seba-pozornosti je pomerne nízka výpočetná komplexita vzhľadom na počet vrstiev v modeli. Ďalšou výhodou je fakt, že transformery nemajú problém s učením sa vzdialených závislostí (long-range dependencies). Vzdialené závislosti sú závislosti medzi dátami, ktoré prekračujú jednotlivé domény. Napr. ak vo vetách „Vyrastal som vo Nemecku. Hovorím po nemecky.“ ak chceme predpovedať posledné slovo druhej vety potrebujeme na to informáciu z predošlej vety (inej domény). S týmto zvyknú mať rekurentné modely problém. [44] Posledným vylepšením je schopnosť paralerizácie učenia. Paralerizácia je niečo, čoho rekurentné modely neboli schopné, čo malo veľký dopad na škálu možností učenia. Transformerom sa podarilo tento problém vyriešiť, čo umožnilo učiť modely na veľkých dátových sádach s výbornými výsledkami. [44] Vďaka tomuto mohli vzniknúť komprehensívne jazykové modely ako je BERT, GPT-3 alebo GPT-4. BERT bol trénovaný na 3.3 miliardách slov. [5]

2.4.1 Transformery v oblasti superrozlišenia

V roku 2021 bol vytvorený transformer, ktorého úlohou bolo klasifikovať obrázky. Tzv. vidiaci transformer (vision transformer - ViT) používa kódér podobný bežnému transformeru, z ktorého sa informácie následne posielajú ďalej na klasifikáciu. [8] Odvtedy bolo vytvorených veľa modelov pre rôzne oblasti počítačového videnia (klasifikácia, segmentácia, a pod.). Všeobecne sa modely transformerov na superrozlišenia skladajú z troch častí

- Extraktor plytkých črt
- Extraktor hlbokých črt
- Prezorkovací blok

Pričom sa architektúra transformeru zvykne uplatňovať práve v extraktore hlbokých črt. [2]



Obr. 2.3: Schéma modelu transformeru s využitím na superrozlíšenie. Zdroj [2] (prekreslené)

Medzi známe modely patrí SwinIR [24], ESRT [28], HAT [4]. Tieto modely sú zamerané na všeobecné superrozlíšenie, t.j. nemusia byť vhodné pre superrozlíšenie tváří. Transformerové modely zamerané na superrozlíšenie tváří nám nie sú známe. Výhodou transformerov je, že nezvyknú byť náročné na GPU. Obzvlášť dobrým príkladom tohoto je ESRT, ktorý zaberá len 4GB grafickej pamäte. [28]

3 Superrozlíšenie tváří

3.1 Unikátne výzvy FSR

Pod pojmom super-rozlíšenie tváří rozumieme aplikáciu superrozlíšenia na vylepšenie kvality ľudských tvarí. Ľudská tvár je vysoko komplexným objekt s mnohými štruktúrami, vlastnosťami, a pod. Komplikovanosť superrozlíšenia tvarí spôsobila, že sa táto problematika rieši ako individuálny problém. Existujú technológie špeciálne vytvorené pre túto záležitosť (URDGN, WaveletSRNet, SPARNet, Super-FAN, atď.). [46] Ciele technológii super-rozlíšenia tvarí možno rozdeliť na 2 kategórie,

- Vylepšenie kvality za účelom estetickým
- Vylepšenia kvality za účelom identifikácie

V prvom prípade sa jedná iba o vizuálne vylepšenie kvality. Nie je tu podstatné udržanie pôvodnosti snímky alebo identifikačných vizuálnych črt. Kľúčová je výsledná kvalita snímky.

Druhý prípad sa, presne naopak, zaoberá zachovaním identifikačných vizuálnych črt. Úlohou je zjednodušenie identifikácie osoby zachytenej na *LR* snímke pomocou super-rozlíšenia. V tomto prípade nie je až tak podstatná estetická kvalita.

Obidva prípady predstavujú rozdielne oblasti výskumu. Táto práca sa zaoberá iba tým druhým.

3.1.1 Kvalitatívne a kvantitatívne meranie kvality

Jedným z hlavných implikácií rozdielov medzi superrozlíšením za účelom identifikácie a superrozlíšením za účelom estetickým je fakt, že vysoké výsledky kvantitatívnych metód merania kvality (PSNR, SSIM, atď.) nemusia zaručiť dostatočné zachovanie vizuálnych identifikačných črt, čo dokáže sťažiť identifikáciu. Na potlačenie tohto efektu sa môžu použiť kvalitatívne metódy hodnotenia kvality. Často zvyknú byť v podobe dotazníka, kde sa respondenti rozhodujú, ktorá snímka sa najviac podobá na SR. Bolo dokázané, že rozdiel medzi kvantitatívnymi a kvalitatívnymi výsledkami je často veľký. [37]

3.1.2 Problémy GAN

Napriek tomu, že GAN dosahujú výborné výsledky v oblasti superrozlíšenia. Nie sú schopné dôveryhodne rekonštruovať obraz. GAN zvyknú dovytvárať neexistujúce detaily, resp. vynechajú iné, podstatné, existujúce detaily. Toto nie je problém pre

superrozlíšenie za účelom estetickým, avšak pre superrozlíšenie za účelom identifikácie to môže byť fatálne. Veľa technológií takéhoto typu je preto pre našu prácu nepoužiteľných. Až extrémnym príkladom tohto fenoménu je technológia PULSE. [32] Tento model je založený na opačnom princípe, ako obvyklé modely super-rozlíšenia. PULSE sa nesnaží naučiť mapovať LR na HR ale naopak, generuje umelo vytvorené HR snímky ľudských tvári a následne im znižuje rozlíšenie, pričom jednotlivé snímky porovnáva s vstupnou LR snímkou. Snímka, ktorá je, po znížení rozlíšenia, najviac podobná vstupnej LR snímke je výstupom tejto technológie. Napriek tomu, že PULSE dokáže generovať \widehat{HR} snímky s vysokou kvalitou, jedná sa o arbitrárne vygenerované snímky, a teda nie je možné aby v nich bola zachovaná pôvodná identita. Zároveň, tieto snímky zvyknú byť veľmi odlišné od požadovaných výsledkov, ako možno vidno na tejto [fotke](#). [37]

3.2 Superrozlíšenia s vopred známymi informáciami

Na uľahčenie práce so snímkami komplexných objektov ako sú ľudské tváre, začali vznikáť spôsoby ako zachytávať niektoré charakteristiky a poskytnúť ich technológiám superrozlíšenia. Medzi tieto charakteristiky patria - orientačné body (facial landmarks), tepelné mapy (facial heatmaps), parsovacie mapy (facial parse maps), atribúty tváre (facial attributes), informácia o identite (identity information). Tieto charakteristiky sa zvyknú tiež nazývať vopred známe informácie (piori information). [13]

3.2.1 Orientačné body

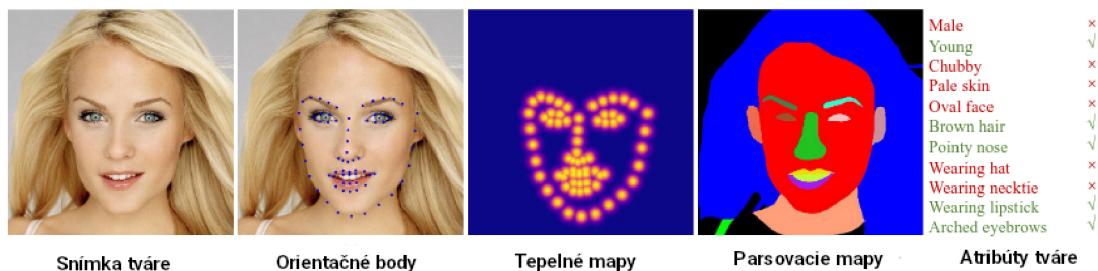
Určujú kľúčové komponenty tváre. Umožňujú technológiám superrozlíšenia lepšie sa orientovať v snímke. Je nutné ich vopred preddefinovať vo fotkách. Existujú databázy ktoré obsahujú veľký počet týchto bodov, napr. Helen [21]

3.2.2 Tepelné mapy

Sú založené na orientačných bodoch. Narozdiel od nich neurčujú priamo kde sa kľúčové komponenty nachádzajú ale určujú pravdepodobnosť, že daný bod bude kľúčový komponent. [13]

3.2.3 Parsovacie mapy

Jedná sa segmentačné mapy, ktorých úlohou je oddeliť od seba sémanticky rozdielne časti tváre, ako sú oči, pery, obočie, koža, atď. [13]



Obr. 3.1: Ukážka vopred známych informácií. Zdroj [13] (preložené)

3.2.4 Atribúty tváre

Atribúty, ako sú pohlavie, má/nemá rúž, má/nemá okuliare, atď. Sa môžu využiť v technológiách superrozlíšenia. Tieto atribúty totižto môžu byť problematické na určenie z LR snímky, a teda sieť nemusí byť schopná ich rekonštrukcie v \widehat{HR} snímke. Vďaka atribútom tváre je možné priamo technológií povedať, čo má v snímke očakávať, čím môžeme zvýšiť presnosť s pôvodnou HR snímkou. Na zjednodušenie sa zvyknú brať v úvahu binárne atribúty. [13]

3.2.5 Informácia o identite

Každá ľudská tvár zodpovedá unikátnej identite. Tento druh informácie je možný zaznamenať a použiť na zlepšenie rekonštrukcie HR snímky. Zrekonštruovaná \widehat{HR} snímka by mala zodpovedať tej istej identite ako pôvodná HR snímka. Informácia o identite sa ukázala, ako veľmi užitočný spôsob zlepšenia výsledkov technológií superrozlíšenia. [13] [47] [11]

3.3 Superrozlíšenie so zachovaním identity

Pod pojmom superrozlíšenie so zachovaním identity rozumieme superrozlíšenie za pomoci znalosti informácie o identite. V posledných rokoch sa tejto oblasti venovalo viac pozornosti, keďže sa ukázalo, že je schopná priniesť výborných výsledkov. [47] [11] Jedná sa zároveň o veľmi zaujímavú oblasť z pohľadu tejto práce, keďže zachovanie identity je kľúčové pre následnú identifikáciu človeka na snímke.

Zachovanie identity môže byť dosiahnuté dvomi spôsobmi, a t.j. zachovanie identity pomocou rozpoznávania tváre a zachovanie identity párovými metódami.

3.3.1 Zachovanie identity pomocou rozpoznávania tváre

Tento systém používa sieť na rozpoznávanie identity tváří aby bolo možné určiť identitu snímok. Na základe určenia identity HR a \widehat{HR} snímky je následne možné vypočítať rozdiel v identite medzi HR a \widehat{HR} . Na základe tohto poznatku môžeme definovať tzv. identitovú stratu (identity loss),

$$\mathcal{L}_{\text{identity}}(\widehat{HR}, HR) = \|FR(\widehat{HR}) - FR(HR)\|_k, \quad (3.1)$$

kde FR reprezentuje funkciu siete na rozpoznávanie identity tváří, k je buď 1 alebo 2, čo závisí od technológie. [13]

Aplikácia tejto straty závisí od individuálnej technológie. Jednou z najjednoduchších je SiCNN (Super-identity convolutional neural network). Tento model používa 2 CNN. Prvá prevzorkuje snímku na vyššie rozlíšenie, zatiaľ čo druhá aplikuje metódu zachovania identity na prevzorkovanú snímku, aby jej znovu zlepšila kvalitu na základe identity. [47]

3.3.2 Zachovanie identity párovými metódami

Najväčším problémom technológií založených na zachovaní identity pomocou rozpoznávania tváre je práve sieť na rozpoznávanie identity tváří. Tieto siete potrebujú veľmi dobre označené dátové sady. To znamená, že každá snímka má k sebe priradenú presnú identitu. Tvorba takýchto dátových sád je veľmi obtiažna. [13]

Na riešenie tohto problému sa môžu použiť tzv. slabo označené párové dátové sady. Toto označenie hovorí len o tom či snímky z daného páru prislúchajú k rovnakej identite, bez toho aby bolo nutné poznať presnú identitu tváří v snímkach. [11]

4 Popis experimentu

4.1 Dátové sady

Jednou z kľúčových častí strojového učenia sú dáta. Modely strojového učenia potrebujú dáta na to, aby boli schopné si nastaviť vnútorné parametre, tak aby dokázali pracovať s čo najväčšou presnosťou. Rovnako to funguje pri superrozlíšení tváří založenom na neurálnych sieťach.

Je empiricky dokázané, že technológie natréňované na snímkach ľudských tváří dosahujú lepšie výsledky pri superrozlíšení tváří ako keď sú natréňované na bežných snímkach. [46] Preto budeme výber dátových sád zameriavať výhradne na sady, ktoré obsahujú výhradne ľudské tváre. Dátové sady obsahujú veľa snímkov ľudských tváří. Niektoré obsahujú len *HR* snímky a *LR* si musíme vyrobiť sami, zatiaľčo iné obsahujú páry, resp. skupiny *HR* a *LR* snímkov. V nasledujúcej tabuľke môžeme vidieť príklady niekoľkých takýchto dátových sád.

Tab. 4.1: Ukážka niektorých dátových sád.

Meno dátovej sady	počet snímkov
CelebA [27]	202 599
VGGFace [36]	3 310 000
Helen [21]	2 330
MLFDB [37]	14 800
FFHQ [14]	70 000

V našom experimente sme si zvolili dátovú sadu CelebA.

4.1.1 CelebA

CelebA je veľmi známa a osvedčená dátová sada. Obsahuje viac ako 200 tisíc snímkov celebrit. Bola vytvorená pokrokom v metóde predpovedania tvárových atribút. Obsahuje výhradne *HR* snímky. CelebA bola použitá vo veľkom počte štúdií v oblasti superrozlíšenia. Je možné ju považovať za jednu z najpopulárnejších dátových sád. [27]

Pri tvorbe *LR* snímkov sme použili 4 druhy interpolácií – bikubická, bilinéarna, najbližších susedov, lanczos. Na jednotlivé snímky sme použili vždy len jednu interpoláciu, zvolenú náhodným výberom. T.j. každá interpolácia mala 20% šancu, že bude aplikovaná. Tento postup sme zvolili aby sa zamedzilo fenoménu, keď sa sieť naučí pracovať len so snímkami vytvorenými len jedným druhom interpolácie.

Tento fenomén spôsobí, že sieť bude dosahovať horšie výsledky pri iných formách degradácie. [37]

4.1.2 MLFDB

MLFDB je relatívne nová dátová sada, zameraná na superrozlíšenie tvárí s viacerými snímkami. Bola vytvorená z viac ako 300 videí na stránke YouTube a obsahuje viac ako 14 tisíc fotiek. Tieto fotky prislúchajú približne 6 000 rôznym osobám. Dátová sada bola vytvorená pomocou technológií na rozpoznávanie tvárí. Fotky sú rozdelené na osmice. Pričom 1 snímka je *HR* a zvyšných 7 je *LR*. MLFDB otvára dvere do neprebádané oblasti superrozlíšenia tvárí s viacerými snímkami. [37]

4.2 Metódy merania kvality

4.2.1 PSNR

Vrcholový pomer signálu k šumu (Peak signal-to-noise ratio – PSNR) je jedna z najznámejších metrik na porovnávanie kvality snímok. Je založená na priemernej kvadratickej chybe (MSE).

$$\text{MSE}(\widehat{HR}, HR) = \frac{1}{hwc} \sum_{i,j,k} (\widehat{HR}_{i,j,k} - HR_{i,j,k})^2, \quad (4.1)$$

$$\text{PSNR}(\widehat{HR}, HR) = 10 \log_{10} \left(\frac{MAX^2}{MSE} \right), \quad (4.2)$$

kde MAX reprezentuje najvyššiu možnú hodnotu pixela (255 pre 8-bitové snímky), zvyšná notácia je rovnaká ako v 2.2.

Napriek svojej popularite naprieč mnohými oblasťami je veľmi dobre známe, že PSNR nezvykne dobre kolerovať so subjektívnym hodnotením kvality človekom. [13]

4.2.2 SSIM

Meranie indexu štruktúrálnej podobnosti (structural similarity index measure – SSIM) sa, na rozdiel od PSNR, zameriava na porovnávanie štruktúrálnej podobnosti. SSIM meria tri oblasti fotky - jas, kontrast a štruktúra. SSIM možno matematicky zapísať ako,

$$\text{SSIM}(\widehat{HR}, HR) = L(\widehat{HR}, HR) * C(\widehat{HR}, HR) * S(\widehat{HR}, HR), \quad (4.3)$$

kde $L(\widehat{HR}, HR)$, $C(\widehat{HR}, HR)$, $S(\widehat{HR}, HR)$ sú jas, kontrast a štruktúra snímok nasledovne. Tieto časti možno zapísať ako

$$L(\widehat{HR}, HR) = \frac{2\mu_{\widehat{HR}}\mu_{HR} + c_1}{\mu_{\widehat{HR}}^2 + \mu_{HR}^2 + c_1}, \quad (4.4)$$

$$C(\widehat{HR}, HR) = \frac{2\sigma_{\widehat{HR}}\sigma_{HR} + c_2}{\sigma_{\widehat{HR}}^2 + \sigma_{HR}^2 + c_2}, \quad (4.5)$$

$$S(\widehat{HR}, HR) = \frac{\sigma_{\widehat{HR}HR} + \frac{c_2}{2}}{\sigma_{\widehat{HR}}\sigma_{HR} + \frac{c_2}{2}}, \quad (4.6)$$

kde $\mu_{\widehat{HR}}$ je pixel s mediánovou hodnotou \widehat{HR} , μ_{HR} je pixel s mediánovou hodnotou HR , $\sigma_{\widehat{HR}}$ je rozptyl \widehat{HR} , σ_{HR} je rozptyl HR , $\sigma_{\widehat{HR}HR}$ je kovariancia medzi \widehat{HR} a HR a c_1 a c_2 sú stabilizačné premenné. [46]

4.2.3 LPIPS

Naučená percepčná obrázková “patchová” podobnosť (Learned perceptual image patch similarity – LPIPS) je percepčná hodnotiacia metrika. LPIPS porovnáva rozdiely medzi snímkami pomocou extrahovania črt snímok. Na extrakciu sa používajú neurálne siete (napr. VGG, Adam, atď.). Výsledky tejto metriky kolerujú so subjektívnym ľudským hodnotením oveľa viac ako PSNR a SSIM. [48]

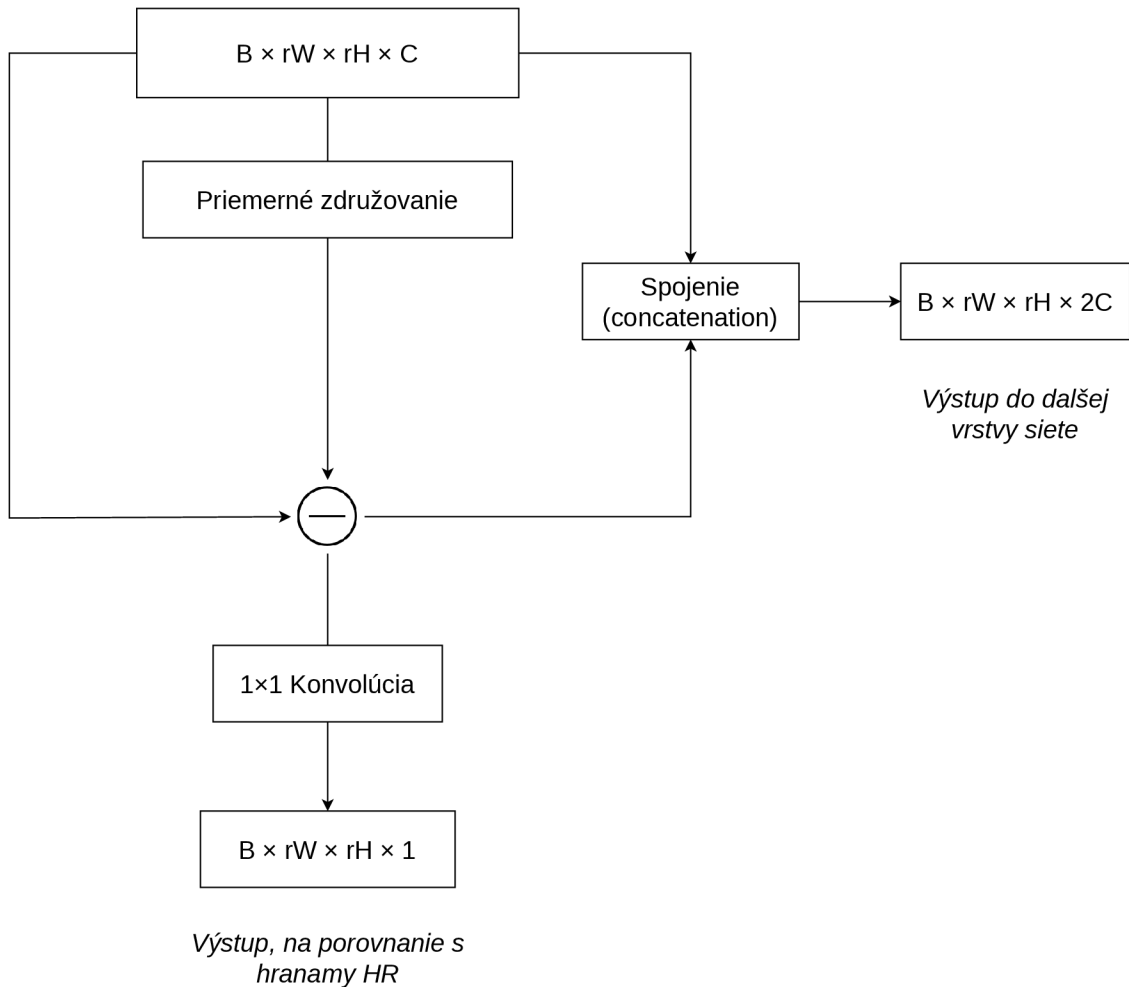
4.3 Zvolene technológie

4.3.1 EIPNet

V roku 2021 bola vytvorená technológia EIPNet. Jedná sa o jednu z najmodernejších technológií v oblasti super-rozlíšenia tváří. EIPNet využíva dvoch vopred známych informácií, a t.j. Informáciu o identite a snímku obrysov tváre. EIPNet vo svojej architektúre využíva tzv. Hranové bloky (edge blocks) na získavanie hrán snímky počas jej prevzorkovania. V sieti je viacero hranových blokov, tzn. hrany snímky sú získavané pri rôznych rozlíšeniach. Výstup týchto blokov je následne porovnávaný s hranami získanými z HR snímky pomocou tzv. cannyho detektoru hrán. Z toho jasne vyplýva, že sieť dokáže merať rozdiely medzi hranami HR a hranami získanými hranovými blokmi, tzn. tento systém prináša tzv. Hranovú stratovú funkciu, ktorú môžeme zapísať ako,

$$\mathcal{L}_e = \frac{1}{r^2wh} \sum_{i,j}^{rw,rh} (C(HR_{i,j}) - E(I_{i,j})) , \quad (4.7)$$

kde I reprezentuje snímku prevzorkovanú o faktor r (nemusí sa jednať o \widehat{HR} , keďže hranových blokov je v sieti niekoľko a všetky okrem posledného majú na vstupe len čiastočne prevzorkovanú snímku), C je cannyho detektor hrán a E je hranový blok. [17]



Obr. 4.1: Nákres hranového bloku. Zdroj [17] (prekreslené)

Zachovanie identity

Okrem Zachovania obrysov sa sieť snaží zachovať aj identity človeka na fotke. Tu sa uplatňuje metóda zachovania identity pomocou rozpoznávania tváre. Konkrétne je tu použitá sieť Inception1, ktorá bola predtrénovaná na dátovej sade VGGFace. [17]

Svetelno-chrominančná strata

EIPNet zároveň počíta tzv. Svetelno-chrominančnú stratu. Autori vychádzajú z ideí, že pri reprezentácii snímku v YUV farebnom modeli je možné dosiahnuť lepšiu percepčnú kvalitu v superrozlíšení, ako pri použití bežného modelu RGB. Technológia si prevedie obrázku z RGB do YUV modelu pomocou konverznej matice a následne vypočíta priemernú kvadratickú chybu v YUV. [17]

Celkovo má stratová funkcia technológie EIPNet tvar,

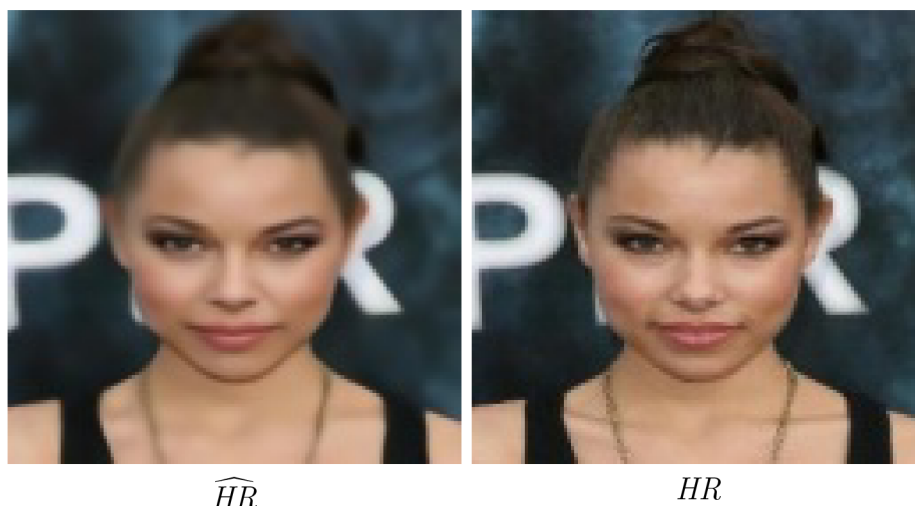
$$\mathcal{L}_{\text{celková}} = \mathcal{L}_2 + \gamma\mathcal{L}_e + \mathcal{L}_{\text{lc}} + \alpha\mathcal{L}_{\text{id}} + \beta\mathcal{L}_{\text{ad}}, \quad (4.8)$$

kde \mathcal{L}_{lc} je svetelno-chrominančná strata, \mathcal{L}_{id} je strata identity, \mathcal{L}_{ad} je adversariálna strata (jedná sa o GAN) a γ, α, β sú váhové parametre. [17]

Trénovanie siete

EIPNet sme natrénovali na CelebA dátovej sade. Sada bola rozdelená na dve množiny – testovaciu a trénovaciu množinu. Pričom trénovacia bola použitá na trénovanie a testovacia na testovanie siete. Testovacia množina obsahovala okolo 60 000 snímok zatiaľčo trénovacia obsahovala okolo 140 000 snímok (t.j. CelebA bola rozdelená približne v pomere 1 : 2).

Na trénovanie sme použili skript `train.py`. Parametre sme zachovali podľa autorov. Počet epoch bol nastavený na 200. Pri 113. epoche sme však trénovanie zastavili, keďže už nedochádzalo k zlepšovaniu. Pred spustením bolo však nutné ešte sieť upraviť, keďže bola naprogramovaná na $8\times$ superrozlíšenie a my pracujeme s $4\times$. Konkrétne bolo nutné upraviť generátor aby fotku prevzorkoval len 4-krát.



Obr. 4.2: Ukažka fotky z modelu EIPNet

4.3.2 WaveletSRNet

Jedná sa o konvolučnú neurálnu sieť, ktorá využíva metód vlnovej transformácie (Wavelet transform). Konkrétne, pred tým než sa sieť pokúsi predpovedať HR snímku sa najprv pokúsi vygenerovať tzv. vlnové koeficienty na základe ktorých sa pokúsi zrekonštruovať HR snímku. Pomocou vlnovej transformácie sa HR snímka pretvorí na vlnové koeficienty, na ktoré sa sieť snaží dostať z LR snímky. Výhodou tejto metódy je, že sieť je schopná lepšie zachytiť a zrekonštruovať perцепčné črty, čo bežná konvolučná sieť nedokáže. WaveletSRNet sa skladá z troch sietí – vkladacia sieť (embedding net), sieť vlnovej predikcie (wavelet prediction net) a rekonštrukčná sieť (reconstruction net). Vkladacia sieť má za úlohu reprezentovať vstupnú snímku ako prízankové mapy. Tie putujú ďalej do siete vlnovej predikcie kde budú predikované vlnové koeficienty. V rekonštrukčnej sieti sa vlnové koeficienty prerobia na \widehat{HR} . Model používa tzv. vlnovú stratu na porovnávanie odchýlky vygenerovaných vlnových koeficientov s požadovanými. [12]

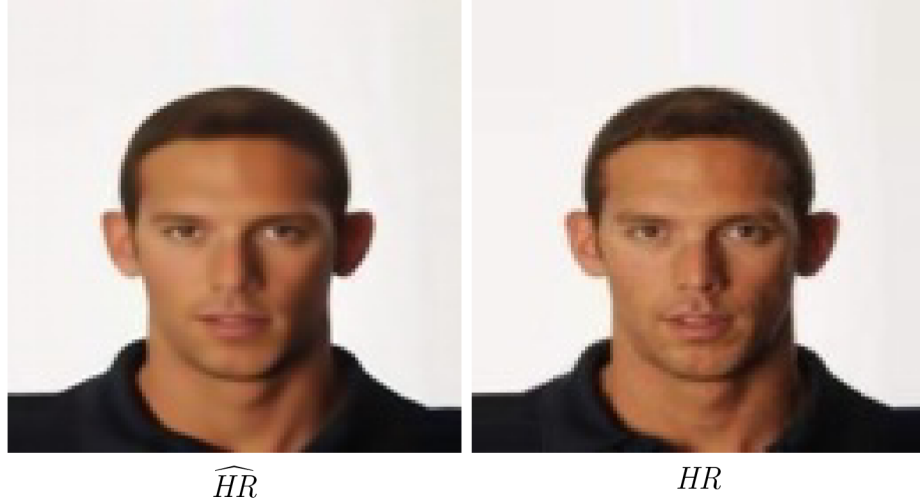
Celkovú stratovú funkciu môžeme zapísať ako,

$$\mathcal{L}_{\text{celkova}} = \mathcal{L}_{\text{vlnova}} + \alpha \mathcal{L}_{\text{texturova}} + \beta \mathcal{L}_2, \quad (4.9)$$

kde $\mathcal{L}_{\text{texturova}}$ je tzv. texturová strata, ktorej úlohou je zabrániť vysokofrekvenčným vlnovým koeficientom aby konvergovali k nule, čím zabraňuje degradácii kvality snímky a α, β sú váhové parametre. \mathcal{L}_2 je taktiež použitá na zachovanie rovnováhy medzi hladkosťou výsledkov a zachovaním podstatných črt. [12]

Tréovanie siete

Sieť sme tréovali na dátovej sade CelebA, ktorú sme rozdelili na testovaciu a tréovaciu množinu. Testovacia množina obsahovala približne 35 000 snímok, zatiaľčo tréovacia obsahovala zvyšných približne 167 000 snímok. Sieť sme tréovali po dobu 200 epoch, zvyšné nastavenia sme nechali rovnako ako v pôvodnej štúdií. Technológia je schopná $4\times$, $8\times$ a $16\times$ superrozlíšenia. Pracovali sme so $4\times$.



Obr. 4.3: Ukažka fotky z modelu WaveletSRNet

4.3.3 NLSN

Nelokálna riedka sieť (Non-local sparse network – NLSN) je model z roku 2021. Je založený na princípoch nelokálnej riedkej pozornosti. Nelokálna pozornosť je špeciálny druh sebaopozornosti, ktorá, zjednodušene povedané, počíta podobnostné skóre medzi párom pozícií na obrázku. Touto metódou je možné lepšie zachytávať vzdialené závislosti. NLSN priniesla nelokálnu riedku pozornosť. Táto forma pozornosti sa snaží potláčať obmedzený rozsah medzi porovnávanými pozíciami. [31]

Ako stratovú funkciu sme zvolili

$$\mathcal{L}_{\text{celková}} = 2 * \mathcal{L}_1 + 0.6 * \mathcal{L}_{\text{ssim}} + 0.2 * \mathcal{L}_{\text{vgg}}, \quad (4.10)$$

kde $\mathcal{L}_{\text{ssim}}$ je stratová funkcia založená na metrike ssim a \mathcal{L}_{vgg} je percepčná stratová funkcia založená na sieti VGG.

$\mathcal{L}_{\text{ssim}}$ má tvar

$$\mathcal{L}_{\text{ssim}} = 1 - SSIM(\widehat{HR}, HR) \quad (4.11)$$

kde $SSIM(\widehat{HR}, HR)$ je metrika SSIM z \widehat{HR} a HR snímok.

Úlohou $\mathcal{L}_{\text{ssim}}$ je, okrem udržania štruktúry, aj potláčanie šumu, ktorý spôsobuje \mathcal{L}_{vgg} . Šum sa nám bohužiaľ nepodarilo kompletne odstrániť. Skúšali sme aj tzv. stratu celkového rozptylu (total variation loss) [46], avšak tá príliš zhoršovala kvalitu snímok. Napriek čiastočnému zašumeniu si podľa nášho názoru zachovávajú dobrú mieru vizuálnej kvality.

Trénovanie siete

Sieť sme trénovali na dátovej sade CelebA. Trénovali sme na približne 130 tisíc snímkach, zatiaľ čo testovanie prebehlo na približne 70 tisíc snímkach. Sieť sme trénovali po dobu 200 epoch.



Obr. 4.4: Ukažka fotky z modelu NLSN

4.3.4 HAT

Hybridný poznostný transformer (Hybrit attention transformer – HAT) je technológia z roku 2023. Jedná sa o jeden z najmodernejších modelov v oblasti vidia-cich transformerov. Je založený na rovnakej schéme ako 2.3. Pričom na extrakciu plytkých črt sa používa bežná konvolúcia. Extraktor hlbokých črt sa skladá z tzv. reziduálnych hybridných pozorostných grúp (Residual Hybrid Attention Groups – RHAG). Tieto bloky obsahujú viacero hybridných pozorostných blokov (Hybrid Attention Block – HAB). HAB kombinuje kanálovú pozorost (channel attention) so seba-pozorostou. Autori vychádzajú z predpokladu, že kombinácia konvolúcií (v kanálovej pozorosti) a seba-pozorosti dokážu vylepšiť výsledky modelu. Na záver, sa na prevzorkovanie používa vrstva Sub-pixel. [28]

Ako stratovú funkciu sieť používa \mathcal{L}_1 , táto funkcia nebola dostatočné pre naše účely. Na superrozlišení tvárí sme použili nasledujúcu stratovú funkciu

$$\mathcal{L}_{\text{celková}} = 2 * \mathcal{L}_1 + 0.8 * \mathcal{L}_{\text{ssim}} + 0.2 * \mathcal{L}_{\text{vgg}} + 0.5 * \mathcal{L}_{\text{sobel}}, \quad (4.12)$$

kde $\mathcal{L}_{\text{ssim}}$ je stratová funkcia založená na metrike ssim, \mathcal{L}_{vgg} je precepčná strata a $\mathcal{L}_{\text{sobel}}$ je hranová strata.

Hranová strata

Ako hranovú stratu používame Sobelovú detekciu hrán (Sobel edge detection). [42] Pracujeme so snímkami s veľkosťami 128×128 pixelov. Tieto snímky sú natolko malé, že Sobelov operátor má problémy správne detekovať mnohé črty. Preto pre tým než aplikujeme na snímky Sobelov operátor ich prevzorkujeme na veľkosť 1024×1024 pomocou bikúbickej interpolácie.



Obr. 4.5: Ukážka rozdielu zachytenia obrysov obrázku po aplikácii sobelovho operátora pri rozdielnych veľkostiach

Nevýhodou tohto prístupu je, že práca so snímkami veľkého rozlíšenia je pamäťovo náročná a zároveň spomaľuje trénovanie siete. Zároveň prevzorkovanie snímku spôsobí, že snímky už nebudú identické, keďže bikúbická interpolácia si dopočíta pixely, ktoré tam predtým neboli. V snímkach však nebol veľký rozdiel takže to nepovažujeme za veľký problém. Aplikáciou tejto metódy detekcie hrán sa nám podarilo zvýšiť PSNR fotiek v priemere o pár desiatín dB.

Podobne ako v predošlom modeli aj tu snímky trpia čiastočným zašumením, ktoré sa nám úplne nepodarilo odstrániť ani pomocou SSIM. Použitie straty celkového rozptylu sa aj v tomto prípade ukázalo ako neúspešné. Našťastie šum ani v tomto prípade nie je veľký a nemyslíme si, že bude mať veľký vplyv na vizuálnu kvalitu snímok z pohľadu identifikácie.

HAT sa ukázal ako model schopný reštaurovať snímky ľudských tvárí s dobrou kvalitou. To svedčí o potenciále architektúr typu transformer v oblasti superrozlíšenia ľudských tvárí. Model zaostáva za ostatnými iba v oblasti zachovania odtieňov tváre.

Trénovanie siete

Model sa (aj bez našej metódy detekcie hrán) trénoval pomaly. Trénovali sme ho len 10 epoch, no to bolo dostačujúce na dosiahnutie dobrých výsledkov. Trénovali

sme na dátovej sade CelebA. Model sme trénovali na približne 130 tisíc obrázkoch, testovanie následne prebehlo na približne 70 tisíc obrázkoch.



\widehat{HR}

HR

Obr. 4.6: Ukažka fotky z modelu HAT

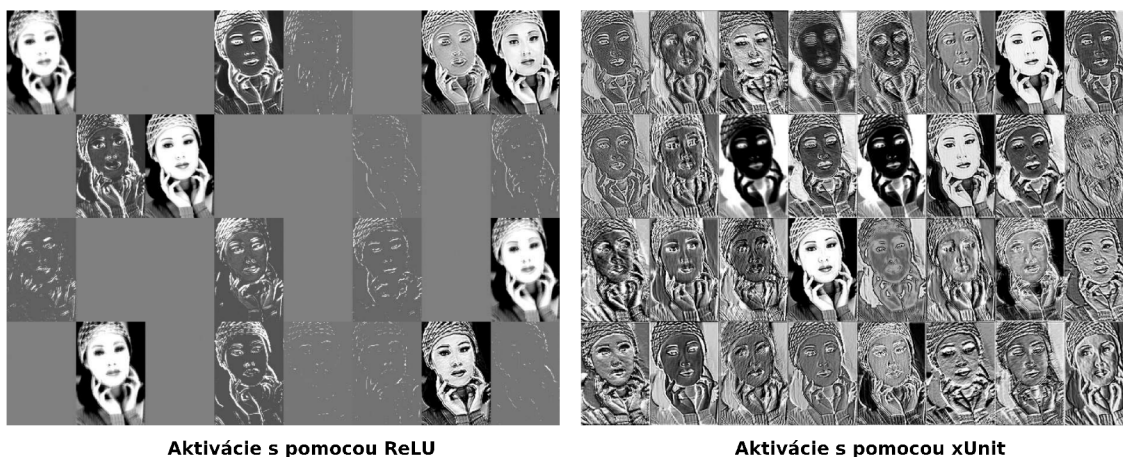
4.3.5 SRDD*

Superrozlíšenie pomocou hlbokého slovníku (super-resolution using deep dictionary - SRDD) je pomerne nová metóda. Na rozdiel od bežných metód super rozlíšenia, je založená na tzv. metóde sparse-coding. Sieť sa učí vytvárať tzv. slovník, ktorý sa používa na premenu LR snímky na \widehat{HR} . Na rozdiel od iných sparse-coding metód sa v tomto prípade slovník vytvára z náhodného šumu. [29]

Pri práci so sieťou sme pozmenili pôvodnú architektúru za cieľom dosiahnutia lepších výsledkov, preto budeme našu sieť nazývať SRDD*. V nasledujúcich častiach prejdeme niektoré zo zmien.

Aktivačná funkcia xUnit

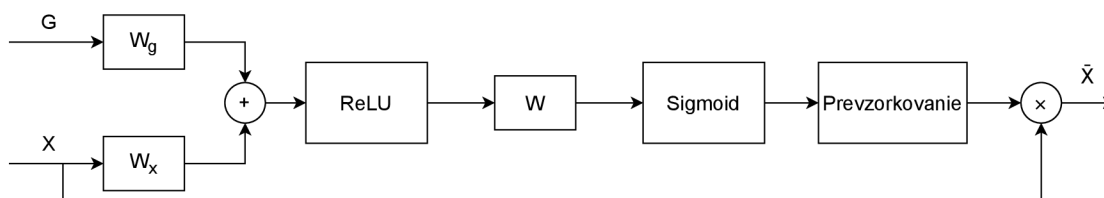
Aktivačná funkcia xUnit je nový druh trénovateľnej aktivačnej funkcie, ktorá dokáže zachytávať viac komplexných črt ako bežné aktivačné funkcie. [18] Experimentáciou sa ukázalo, že SRDD* dosahuje najlepšie výsledky ak xUnit použijeme v extraktore hlbokých črt a pri rekonštrukcii \widehat{HR} snímky.



Obr. 4.7: Porovnanie aktivácií na sieti SRCNN s pomocou ReLU a s pomocou xUnit. Zdroj [18] (preložené)

AttentionUNet++

SRDD používa upravenú UNet++ ako extraktor hlbokých črt. Rozhodli sme sa ju prerobiť na UNet++, ktorá používa pozornosť. Pozornostná UNet++ sa ukázal ako výhodný spôsob segmentácie obrazu hlavne na medicínske účely. [35] Architektúru UNet++ sme ponechali rovnakú ako v pôvodnom SRDD, len sme pridali pozornostné brány na reziduálne prepojenia, rovnako ako v [35].



Obr. 4.8: Schéma pozornostnej brány. Zdroj [35] (prekreslené)

Brána má na vstupe dva signály – X a G. X je vstup z reziduálneho prepojenia, zatiaľ čo G je vstup z uzlu, ktorý sa nachádza „pod“ bránou. Na vstupy sa aplikujú lineárne vrstvy W_x a W_g . Tie sa následne sčítajú a na ich súčet sa aplikuje aktivačná funkcia ReLU. Na výstup z ReLU sa použije ďalšia lineárna vrstva a jej výstup sa ďalej pošle do funkcie Sigmoid. Výstup zo funkcie Sigmoid je následne prevzorkovaný na potrebnú veľkosť a na záver je spojený (concatenated) so signálom X. [35]

Použili sme nasledovnú stratovú funkciu

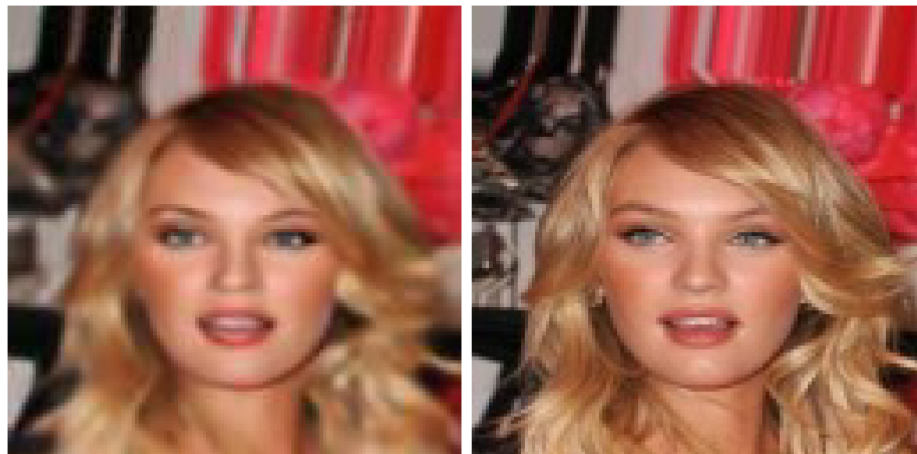
$$\mathcal{L}_{\text{celková}} = 2 * \mathcal{L}_1 + 0.8 * \mathcal{L}_{\text{ssim}} + 0.5 * \mathcal{L}_{\text{LoG}}, \quad (4.13)$$

kde \mathcal{L}_{LoG} je hranová strata založená na LoG (laplacian of gaussian) operátore. [1] LoG operátor sme použili, keďže bežný Sobelov operátor zachytával zbytočne veľa detailov. Detekciu hrán sme aplikovali rovnako ako v 4.3.4.

Na rozdiel od iných sietí sme v tomto prípade nepoužili VGG stratovú funkciu. VGG opäť vytvárala pomerne veľký šum, čo degradovalo kvalitu fotky. Experimentálne sme zistili, že SRDD* dokáže vytvárať fotky s relatívne dobrou kvalitou ja bez VGG stratovej funkcie. Preto sme sa ju rozhodli odstrániť. Nevýhodou tohto prístupu je, že bez VGG stratovej funkcie sa niektoré fotky javia čiastočne rozmazane (viď 4.9). Zároveň, kvôli nepoužitiu percepčnej straty, dosahujú fotky horší LPIPS. Nevidíme v tom však veľký problém, keďže snímky, napriek tomu, vyzerajú podobne z pohľadu identifikácie.

Trénovanie siete

Sieť sme trénovali na rovnakých parametroch ako predošlé siete. Množinu CelebA sme rozdelili na 130 tisíc trénovacích snímok, zvyšok (okolo 70 tisíc) sme použili na testovanie. Trénovali sme po dobu 200 epoch. Trénovanie trvalo približne 4 dni.

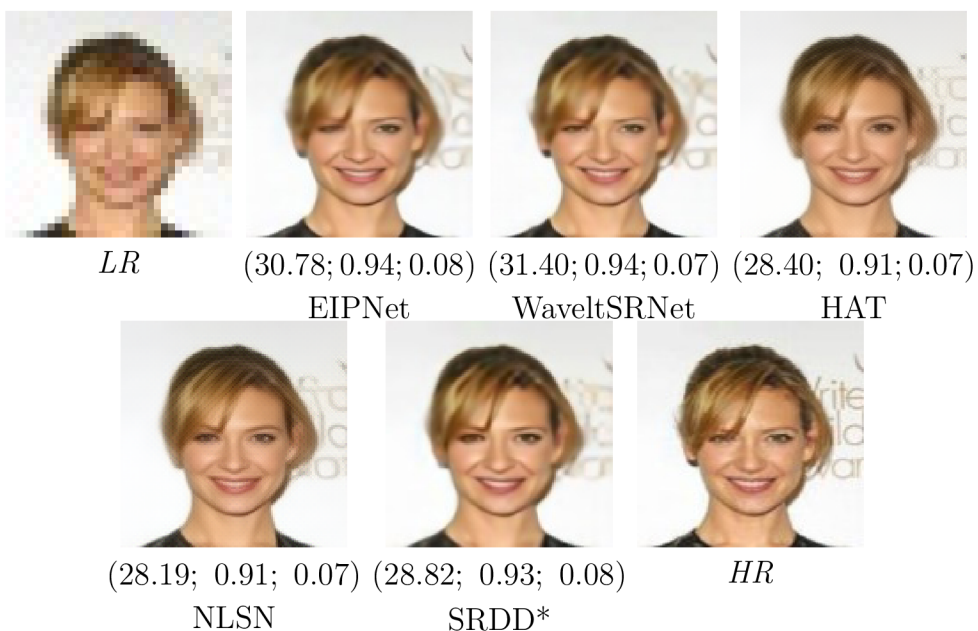


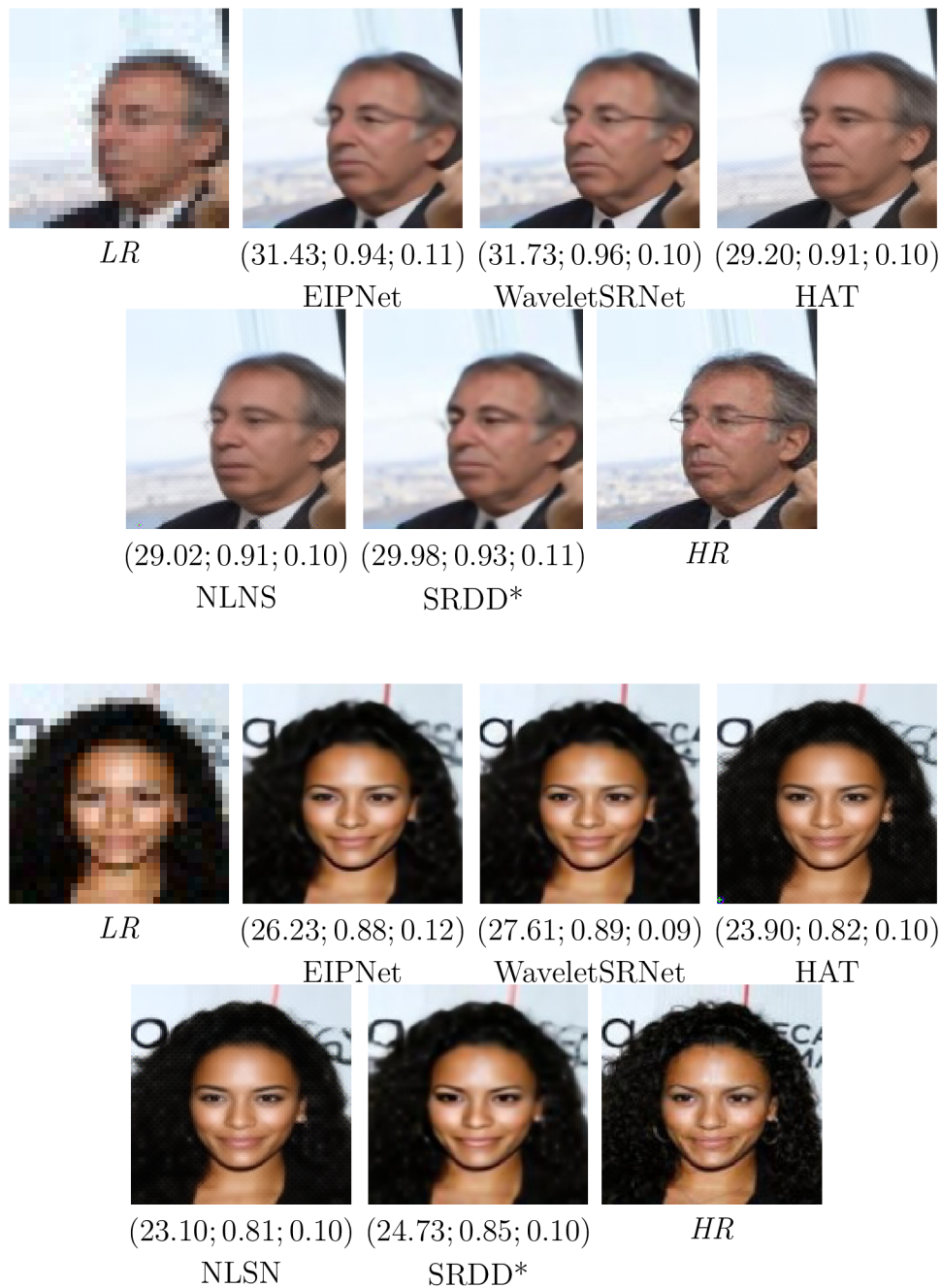
\widehat{HR} HR
Obr. 4.9: Ukážka výstupu modelu SRDD*

5 Výsledky experimentu

Tab. 5.1: Porovnanie výsledkov na dátovej sade CelebA pri 4× superrozlíšení. Najlepšie hodnoty sú **červenou**.

CelebA			
Model	PSNR [dB]	SSIM	LPIPS
EIPNet	28.87	0.92	0.119
WaveletSRNet	29.42	0.93	0.109
NLSN	27.34	0.90	0.099
HAT	27.42	0.90	0.097
SRDD*	28.20	0.91	0.113





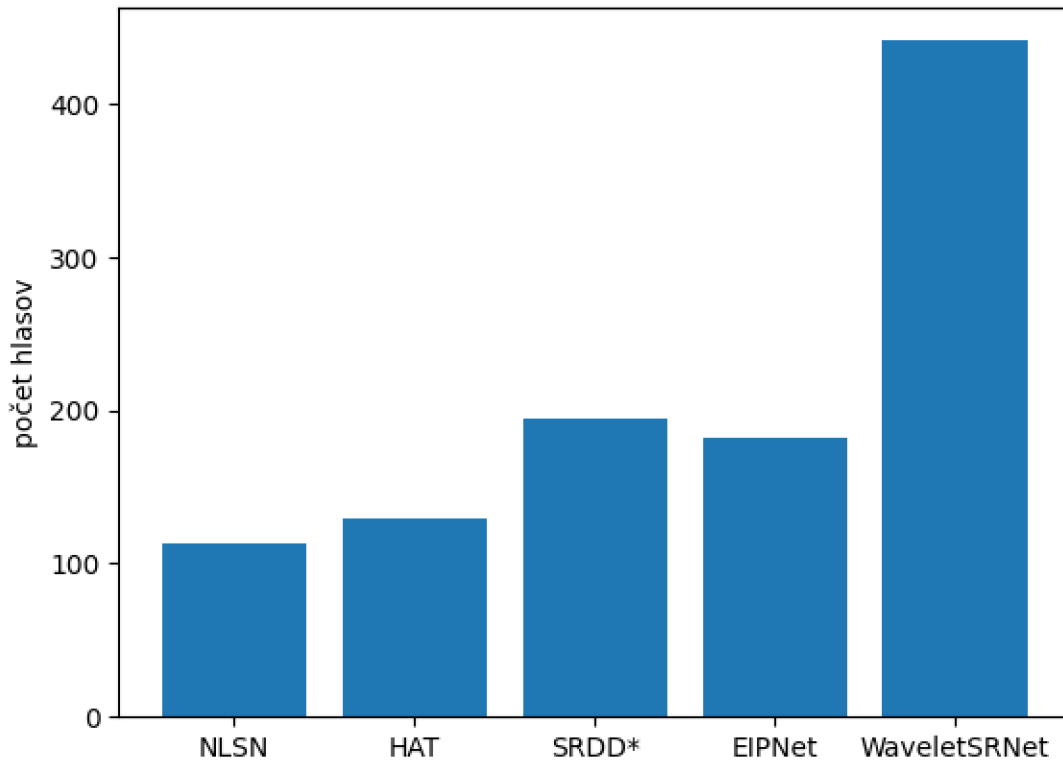
Obr. 5.1: Kvalitatívne porovnanie výsledkov. V zátvorkách sú PSNR [dB], SSIM a LPIPS nasledovne.

Trénovanie a testovanie technológií prebehlo na NVIDIA Tesla v100s grafickej karte. Kvantitatívne výsledky sme merali pomocou troch metrik – PSNR, SSIM a LPIPS.

Kvantitatívne výsledky sa ukázali ako uspokojivé. Všetky siete dosiahli výsledkov podobných s momentálnym stavom vedy a techniky superrozlišenia. NLSN a HAT dosiahli výborných výsledkov v metrike LPIPS. Čo naznačuje výbornú vizuálnu kvalitu. Tieto siete zároveň dosiahli horších výsledkov v metrikách PSNR a SSIM.

Kvalitatívne hodnotenie vyzerá na prvý pohľad ako dostatočné pri väčšine prípadov, pre lepšie hodnotenie sme zostrojili dotazník 5.2. Môžeme predpokladať, že kvalitatívne hodnotenie snímok bude mať väčší význam ako kvantitatívne, vzhľadom na zameranie práce. [37]

5.1 Dotazník



Obr. 5.2: Výsledky dotazníku

Dotazník sme vytvorili náhodným výberom snímok z testovacej časti množiny CelebA. Respondenti v dotazníku porovnávali vizuálnu podobnosť snímok z jednotlivých technológií s pôvodnou \widehat{HR} snímkou a následne zvolili tú, ktorá im prišla najpodobnejšia. Snímok sme zvolili iba 10, keďže sme predpokladali, že veľký počet snímok by mohol respondentov odradiť. Dotazník bol vyplnený 106 respondentmi (t.j. pri 10 snímkoch bolo v dotazníku dokopy 1060 hlasov).

Dotazníku jasne dominovala WaveletSRNet (441 hlasov). Za ňou sa umiestnila naša technológia SRDD* (194 hlasov). EIPNet mala o 12 hlasov menej (182 hlasov). Na posledných dvoch miestach sa umiestnili HAT (130 hlasov) a NLSN (113 hlasov). Predpokladáme, že HAT a NLSN skončili na posledných miestach, keďže nedokázali dobre zachytávať odtiene tváre. Zaujímavé je, že na posledných 2 miestach

sa umiestnili technológie s najlepším LPIPS. čo ukazuje, že aj metriky snažiace sa napodobniť ľudské vnímanie majú svoje nedostatky. K tomuto mohol prispieť aj faktor, že snímky z HAT a NLSN boli čiastočne zašumené.

Experiment ukázal účinnosť vlnovej transformácie v oblasti superrozlíšenia. Na prvých miestach sa umiestnili technológie, ktoré boli schopné zachytávať farby a od-tiene pokožky. Preto považujeme tento faktor za veľmi podstatný pri tvorbe nových technológií superrozlíšenia tváří.

Použitie pozornosti v SRDD* prinieslo len o niečo lepšie výsledky ako pôvodný model. Aplikácia pozornosti na superrozlíšenie je preto problematická. Napriek tomu náš model prekonal všetky ostatné, s výnimkou WaveletSRNet. Predpokladáme, že pre lepšie výsledky by mohli byť použité komplexnejšie pozornostné bloky.

Taktiež je nutné povedať, že práca sa nevenuje tzv. skutočnému superrozlíšeniu (real super-resolution). Pod týmto pojmom sa rozumie superrozlíšenie z LR snímok získaných skutočnou degradáciou kvality (napr. kamera s nízkym rozlíšením, šum, nevhodné osvetlenie, zlé počasie, a pod.), na rozdiel od bežne používaných interpolácií. Na túto problematiku bohužiaľ neexistujú dostatočné dátové sady. Toto môže mať negatívny efekt na aplikovateľnosť daných sietí v praxi. [3]

Záver

V rámci našej práce sme sa venovali technológiám superrozlíšenie a následne sme natrénovali a porovnali 5 rôznych neurálnych sietí. Jednou z týchto sietí bola nami upravená verzia. Všetky siete boli natrénované a podrobne sme ich kvantitatívne aj kvalitatívne porovnali. Výsledky kvantitatívnej analýzy ukázali, že všetky siete dosahujú konkurencieschopné výsledky v porovnaní s modernými štandardmi v tejto oblasti. Teda všetky siete sa ukázali ako účinné metódy pre superrozlíšenie.

Objektívne aj subjektívne výsledky ukázali, že siete dosahujú dobré hodnoty percepčnej kvality, čo naznačuje ich potenciál pre identifikáciu. Výsledky sme porovnali pomocou dotazníka, z ktorého vyplynulo, že vlnová transformácia má výnimočný potenciál v oblasti superrozlíšenia tvarí.

Dotazníkové hodnotenie tiež ukázalo, že zachovanie odtieňov a farieb tváre je dôležité pre subjektívne hodnotenie. V tejto oblasti sa najlepšie výsledky dosiahli pomocou siete s názvom WaveletSRNet. Táto sieť dokázala prekonať aj technológie, ktoré dosiahli dobré výsledky z hľadiska objektívnych metrík. Nami upravená sieť dosiahla dobrých výsledkov, porovnateľných z modernými technológiami superrozlíšenia.

Je však potrebné zdôrazniť, že naša práca sa nezaoberala reálnym superrozlíšením (real super-resolution). Toto môže mať významný vplyv na aplikovateľnosť týchto technológií v skutočných prípadoch.

Literatúra

- [1] AHMADI, N.; AKBARIZADEH, G. *Iris Recognition System based on Canny and LoG Edge Detection Methods*. *Journal of Soft Computing and Decision Support Systems*, 2, 2015.
- [2] ALI, A. M.; BENJDIRA, B.; KOUBAA, A.; EL-SHAFI, W.; KHAN, Z.; BOULILA, W. *Vision Transformers in Image Restoration: A Survey*. *Sensors*, 23(5), 2023. ISSN 1424-8220. doi:10.3390/s23052385. Dostupné z: <<https://www.mdpi.com/1424-8220/23/5/2385>>.
- [3] CHEN, H.; HE, X.; QING, L.; WU, Y.; REN, C.; SHERIFF, R. E.; ZHU, C. *Real-world single image super-resolution: A brief review*. *Information Fusion*, 79:124–145, 2022. ISSN 1566-2535. doi:<https://doi.org/10.1016/j.inffus.2021.09.005>. Dostupné z: <<https://www.sciencedirect.com/science/article/pii/S1566253521001792>>.
- [4] CHEN, X.; WANG, X.; ZHOU, J.; QIAO, Y.; DONG, C. *Activating More Pixels in Image Super-Resolution Transformer*. 2023.
- [5] DEVLIN, J.; CHANG, M.; LEE, K.; TOUTANOVA, K. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. *CoRR*, abs/1810.04805, 2018. Dostupné z: <<http://arxiv.org/abs/1810.04805>>.
- [6] DONG, C.; LOY, C. C.; HE, K.; TANG, X. *Image Super-Resolution Using Deep Convolutional Networks*. *CoRR*, abs/1501.00092, 2015. Dostupné z: <<http://arxiv.org/abs/1501.00092>>.
- [7] DONG, C.; LOY, C. C.; TANG, X. *Accelerating the Super-Resolution Convolutional Neural Network*. *CoRR*, abs/1608.00367, 2016. Dostupné z: <<http://arxiv.org/abs/1608.00367>> [cit. 26. november 2022].
- [8] DOSOVITSKIY, A.; BEYER, L.; KOLESNIKOV, A.; WEISSENBORN, D.; ZHAI, X.; UNTERTHINER, T.; DEGHANI, M.; MINDERER, M.; HEIGOLD, G.; GELLY, S.; USZKOREIT, J.; HOULSBY, N. *An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale*. *CoRR*, abs/2010.11929, 2020. Dostupné z: <<https://arxiv.org/abs/2010.11929>>.
- [9] GOODFELLOW, I. J.; POUGET-ABADIE, J.; MIRZA, M.; XU, B.; WARDEFARLEY, D.; OZAIR, S.; COURVILLE, A.; BENGIO, Y. *Generative Adversarial Networks*, 2014. doi:10.48550/ARXIV.1406.2661. Dostupné z: <<https://arxiv.org/abs/1406.2661>> [cit. 26. november 2022].

- [10] GU, J.; WANG, Z.; KUEN, J.; MA, L.; SHAHROUDY, A.; SHUAI, B.; LIU, T.; WANG, X.; WANG, G.; CAI, J.; CHEN, T. *Recent advances in convolutional neural networks*. Pattern Recognition, 77:354–377, 2018. ISSN 0031-3203. doi:<https://doi.org/10.1016/j.patcog.2017.10.013>. Dostupné z: <<https://www.sciencedirect.com/science/article/pii/S0031320317304120>> [cit. 26. november 2022].
- [11] HSU, C.; LIN, C.; SU, W.; CHEUNG, G. *SiGAN: Siamese Generative Adversarial Network for Identity-Preserving Face Hallucination*. IEEE Transactions on Image Processing, 28(12):6225–6236, 2019. doi:10.1109/TIP.2019.2924554.
- [12] HUANG, H.; HE, R.; SUN, Z.; TAN, T. *Wavelet-SRNet: A Wavelet-Based CNN for Multi-scale Face Super Resolution*. In 2017 IEEE International Conference on Computer Vision (ICCV), s. 1698–1706. 2017. doi:10.1109/ICCV.2017.187.
- [13] JIANG, J.; WANG, C.; LIU, X.; MA, J. *Deep Learning-based Face Super-resolution: A Survey*. CoRR, abs/2101.03749, 2021. Dostupné z: <<https://arxiv.org/abs/2101.03749>> [cit. 26. november 2022].
- [14] KARRAS, T.; LAINE, S.; AILA, T. *A Style-Based Generator Architecture for Generative Adversarial Networks*. CoRR, abs/1812.04948, 2018. Dostupné z: <<http://arxiv.org/abs/1812.04948>>.
- [15] KAWULOK, M.; BENECKI, P.; PIECHACZEK, S.; HRYNCZENKO, K.; KOSTRZEWA, D.; NALEPA, J. *Deep Learning for Multiple-Image Super-Resolution*. IEEE Geoscience and Remote Sensing Letters, 17(6):1062–1066, 2020. doi:10.1109/LGRS.2019.2940483.
- [16] KIM, J.; LEE, J. K.; LEE, K. M. *Deeply-Recursive Convolutional Network for Image Super-Resolution*. CoRR, abs/1511.04491, 2015. Dostupné z: <<http://arxiv.org/abs/1511.04491>> [cit. 26. november 2022].
- [17] KIM, J.; LI, G.; YUN, I.; JUNG, C.; KIM, J. *Edge and Identity Preserving Network for Face Super-Resolution*. CoRR, abs/2008.11977, 2020. Dostupné z: <<https://arxiv.org/abs/2008.11977>>.
- [18] KLIGVASSER, I.; ROTT SHAHAM, T.; MICHAELI, T. *xunit: Learning a spatial activation function for efficient image restoration*. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, s. 2433–2442. 2018.
- [19] KO, S.; DAI, B.-R. *Multi-Laplacian GAN with Edge Enhancement for Face Super Resolution*. In 2020 25th International Conference on Pattern Recognition (ICPR), s. 3505–3512. 2021. doi:10.1109/ICPR48806.2021.9412950.

- [20] LAI, W.; HUANG, J.; AHUJA, N.; YANG, M. *Deep Laplacian Pyramid Networks for Fast and Accurate Super-Resolution*. CoRR, abs/1704.03915, 2017. Dostupné z: <<http://arxiv.org/abs/1704.03915>> [cit. 26. november 2022].
- [21] LE, V.; BRANDT, J.; LIN, Z.; BOURDEV, L.; HUANG, T. S. *Interactive facial feature localization*. In European conference on computer vision, s. 679–692. Springer, 2012.
- [22] LEDIG, C.; THEIS, L.; HUSZAR, F.; CABALLERO, J.; AITKEN, A. P.; TEJANI, A.; TOTZ, J.; WANG, Z.; SHI, W. *Photo-Realistic Single Image Super-Resolution Using a Generative Adversarial Network*. CoRR, abs/1609.04802, 2016. Dostupné z: <<http://arxiv.org/abs/1609.04802>> [cit. 26. november 2022].
- [23] LEE, M. *Multi-image Super-resolution via Quality Map Associated Attention Network*, 2022. doi:10.48550/ARXIV.2202.13124. Dostupné z: <<https://arxiv.org/abs/2202.13124>>.
- [24] LIANG, J.; CAO, J.; SUN, G.; ZHANG, K.; VAN GOOL, L.; TIMOFTE, R. *Swinir: Image restoration using swin transformer*. In Proceedings of the IEEE/CVF international conference on computer vision, s. 1833–1844. 2021.
- [25] LIM, B.; SON, S.; KIM, H.; NAH, S.; LEE, K. M. *Enhanced Deep Residual Networks for Single Image Super-Resolution*. CoRR, abs/1707.02921, 2017. Dostupné z: <<http://arxiv.org/abs/1707.02921>> [cit. 26. november 2022].
- [26] LIU, H.; ZHENG, X.; HAN, J.; CHU, Y.; TAO, T. *Survey on GAN-based face hallucination with its model development*. IET Image Processing, 13(14):2662–2672, 2019. doi:<https://doi.org/10.1049/iet-ipr.2018.6545>. Dostupné z: <<https://ietresearch.onlinelibrary.wiley.com/doi/abs/10.1049/iet-ipr.2018.6545>>.
- [27] LIU, Z.; LUO, P.; WANG, X.; TANG, X. *Deep Learning Face Attributes in the Wild*. CoRR, abs/1411.7766, 2014. Dostupné z: <<http://arxiv.org/abs/1411.7766>> [cit. 26. november 2022].
- [28] LU, Z.; LIU, H.; LI, J.; ZHANG, L. *Efficient transformer for single image super-resolution*. arXiv preprint arXiv:2108.11084, 2021.
- [29] MAEDA, S. *Image Super-Resolution with Deep Dictionary*. In Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XIX, s. 464–480. Springer, 2022.

- [30] MEI, Y.; FAN, Y.; ZHOU, Y. *Image Super-Resolution With Non-Local Sparse Attention*. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), s. 3517–3526. 2021.
- [31] MEI, Y.; FAN, Y.; ZHOU, Y. *Image Super-Resolution with Non-Local Sparse Attention*. In 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), s. 3516–3525. 2021. doi:10.1109/CVPR46437.2021.00352.
- [32] MENON, S.; DAMIAN, A.; HU, S.; RAVI, N.; RUDIN, C. *PULSE: Self-Supervised Photo Upsampling via Latent Space Exploration of Generative Models*. CoRR, abs/2003.03808, 2020. Dostupné z: <<https://arxiv.org/abs/2003.03808>> [cit. 26. november 2022].
- [33] MOLINI, A. B.; VALSESIA, D.; FRACASTORO, G.; MAGLI, E. *DeepSUM: Deep Neural Network for Super-Resolution of Unregistered Multitemporal Images*. IEEE Transactions on Geoscience and Remote Sensing, 58(5):3644–3656, 2020. doi:10.1109/tgrs.2019.2959248. Dostupné z: <<https://doi.org/10.1109%2Ftgrs.2019.2959248>>.
- [34] NALEPA, J.; HRYNCZENKO, K.; KAWULOK, M. *Multiple-Image Super-Resolution Using Deep Learning and Statistical Features*. In Andrea Torsello; Luca Rossi; Marcello Pelillo; Battista Biggio; Antonio Robles-Kelly, editoři, Structural, Syntactic, and Statistical Pattern Recognition, s. 261–271. Cham: Springer International Publishing, 2021. ISBN 978-3-030-73973-7.
- [35] OKTAY, O.; SCHLEMPER, J.; FOLGOC, L. L.; LEE, M. C. H.; HEINRICH, M. P.; MISAWA, K.; MORI, K.; MCDONAGH, S. G.; HAMMERLA, N. Y.; KAINZ, B.; GLOCKER, B.; RUECKERT, D. *Attention U-Net: Learning Where to Look for the Pancreas*. CoRR, abs/1804.03999, 2018. Dostupné z: <<http://arxiv.org/abs/1804.03999>>.
- [36] PARKHI, O. M.; VEDALDI, A.; ZISSERMAN, A. *Deep Face Recognition*. In British Machine Vision Conference. 2015.
- [37] RAJNOHA, M.; MEZINA, A.; BURGET, R. *Multi-Frame Labeled Faces Database: Towards Face Super-Resolution from Realistic Video Sequences*. Applied Sciences, 10(20), 2020. ISSN 2076-3417. doi:10.3390/app10207213. Dostupné z: <<https://www.mdpi.com/2076-3417/10/20/7213>> [cit. 26. november 2022].
- [38] SALVETTI, F.; MAZZIA, V.; KHALIQ, A.; CHIABERGE, M. *Multi-Image Super Resolution of Remotely Sensed Images Using Residual Attention Deep*

- Neural Networks*. Remote Sensing, 12(14), 2020. ISSN 2072-4292. doi: 10.3390/rs12142207. Dostupné z: <<https://www.mdpi.com/2072-4292/12/14/2207>> [cit. 26. november 2022].
- [39] SALVETTI, F.; MAZZIA, V.; KHALIQ, A.; CHIABERGE, M. *Multi-Image Super Resolution of Remotely Sensed Images Using Residual Attention Deep Neural Networks*. Remote Sensing, 12(14):2207, 2020. doi:10.3390/rs12142207. Dostupné z: <<https://doi.org/10.3390/rs12142207>>.
- [40] SHI, W.; CABALLERO, J.; HUSZÁR, F.; TOTZ, J.; AITKEN, A. P.; BISHOP, R.; RUECKERT, D.; WANG, Z. *Real-Time Single Image and Video Super-Resolution Using an Efficient Sub-Pixel Convolutional Neural Network*. CoRR, abs/1609.05158, 2016. Dostupné z: <<http://arxiv.org/abs/1609.05158>> [cit. 26. november 2022].
- [41] SIMONYAN, K.; ZISSERMAN, A. *Very deep convolutional networks for large-scale image recognition*. arXiv preprint arXiv:1409.1556, 2014.
- [42] SOBEL, I.; FELDMAN, G. *An Isotropic 3x3 Image Gradient Operator*. 2015. doi:10.13140/RG.2.1.1912.4965.
- [43] TAI, Y.; YANG, J.; LIU, X. *Image Super-Resolution via Deep Recursive Residual Network*. In 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), s. 2790–2798. 2017. doi:10.1109/CVPR.2017.298.
- [44] VASWANI, A.; SHAZEER, N.; PARMAR, N.; USZKOREIT, J.; JONES, L.; GOMEZ, A. N.; KAISER, L.; POLOSUKHIN, I. *Attention Is All You Need*. CoRR, abs/1706.03762, 2017. Dostupné z: <<http://arxiv.org/abs/1706.03762>>.
- [45] WANG, X.; YU, K.; WU, S.; GU, J.; LIU, Y.; DONG, C.; LOY, C. C.; QIAO, Y.; TANG, X. *ESRGAN: Enhanced Super-Resolution Generative Adversarial Networks*. CoRR, abs/1809.00219, 2018. Dostupné z: <<http://arxiv.org/abs/1809.00219>> [cit. 26. november 2022].
- [46] WANG, Z.; CHEN, J.; HOI, S. C. *Deep learning for image super-resolution: A survey*. IEEE transactions on pattern analysis and machine intelligence, 43(10):3365–3387, 2020.
- [47] ZHANG, K.; ZHANG, Z.; CHENG, C.; HSU, W. H.; QIAO, Y.; LIU, W.; ZHANG, T. *Super-Identity Convolutional Neural Network for Face Hallucination*. CoRR, abs/1811.02328, 2018. Dostupné z: <<http://arxiv.org/abs/1811.02328>> [cit. 26. november 2022].

- [48] ZHANG, R.; ISOLA, P.; EFROS, A. A.; SHECHTMAN, E.; WANG, O. *The Unreasonable Effectiveness of Deep Features as a Perceptual Metric*. CoRR, abs/1801.03924, 2018. Dostupné z: <<http://arxiv.org/abs/1801.03924>> [cit. 26. november 2022].

Zoznam symbolov a skratiek

HR	Snímka s vysokým rozlíšením (ground-truth)
LR	Snímka s nízkym rozlíšením
\widehat{HR}	Výstup siete
CNN	Konvolučná neurálna sieť (Convolutional Neural Network)
GAN	Generatívna adversariálna sieť (Generative adversarial network)
ViT	Vidiaci transformer (Vision Transformer)
PSNR	Vrcholový pomer signálu k šumu (Peak signal-to-noise ratio)
SSIM	Meranie indexu štruktúrálnej podobnosti (structural similarity index measure)
LPIPS	Naučená percepčná obrázková “patchová” podobnosť (Learned perceptual image patch similarity)
\mathcal{L}_1	Stredná absolútna chyba (Mean absolute error)
\mathcal{L}_2	Stredná kvadratická chyba (Mean squared error)
ReLU	Lineárny jednotkový usmerňovač (Rectifier linear unit)
xUnit	Aktivačná funkcia xUnit
EIPNet	Sieť zachovávajúca obrys a identitu (Edge and Identity preserving network)
NLSN	Nelokálna riedka sieť (Non-local sparse network)
HAT	Hybridny pozornostný transformer (Hybrid attention transformer)
SRDD	Superrozlíšenia pomocou hlbokého slovníka (super-resolution with deep dictionary)
WaveletSRNet	Vlnová konvolučná neurálna sieť na superrozlíšenie tvarí rôznych škál (A Wavelet-based CNN for Multi-scale Face Super Resolution)