

Evoluční návrh konvolučních neuronových sítí

Michal Piňos

Abstrakt

Cílem této práce je návrh a implementace programu pro automatizovaný návrh konvolučních neuronových sítí (CNN) s využitím evolučních výpočetních technik. Z praktického hlediska tento přístup redukuje potřebu lidského faktoru při tvorbě CNN, a tak eliminuje zdlouhavý a namáhavý proces návrhu. Tato práce využívá speciální formu genetického programování nazývanou kartézské genetické programování, které pro zakódování řešeného problému využívá grafovou reprezentaci. Tato technika umožňuje uživateli parametrizovat proces hledání CNN, a tak se zaměřit na architektury zajímavé z pohledu použitých výpočetních jednotek, přesnosti či počtu parametrů. Navrhovaný přístup byl otestován na standardizované datové sadě CIFAR-10, která je často využívána výzkumníky pro srovnání výkonnosti jejich CNN. Prvotní experimenty ukázaly, že vytvořená implementace (využívající GPU akceleraci) je schopna vytvořit či vylepšit přesnost CNN. Výsledkem experimentů, kdy bylo pro trénování k dispozici pouze několik epoch, byla řešení s přesností 64.5 % a počtem parametrů 146K při využití základních vrstev a řešení s přesností 74.5 % s počtem parametrů 475K při využití reziduálních vrstev. Záměrem těchto experimentů bylo dokázat funkčnost implementovaného programu a *proof-of-concept* navržené metody.

Klíčová slova: Neuroevoluce — Neuronové sítě — Evoluční výpočetní techniky

Příložené materiály: N/A

*xpinos03@stud.fit.vutbr.cz, *Fakulta informačních technologií, Vysoké učení technické v Brně*

1. Úvod

Návrh architektury konvolučních neuronových sítí dnes hraje důležitou roli ve výzkumu umělé inteligence (konkrétně v oboru rozpoznávání obrázků). Návrh nových architektur CNN vyžaduje velké úsilí, mnoho zkušeností a pokusů. Tento proces je tak velmi namáhavý a ne vždy efektivní. Z tohoto důvodu je nutné přemýšlet o nových způsobech tvorby těchto sítí, které by nebyly tak složité a náročné, a přitom poskytovaly výsledky minimálně srovnatelné s ručně navrženými sítěmi. Technika automatizovaného hledání architektur umělých neuronových sítí, nazývaná NAS (Neural Architecture Search) [1], je v dnešní době velmi populární [2, 3] a již dala vzniknout hlubokým neuronovým sítím překonávajícím ty odborně navržené [4, 5, 6].

Tvorba nových architektur CNN rovněž obnáší problém optimalizace navrhované sítě z hlediska počtu parametrů, hyperparametrů (počtu vrstev, skrytých neuronů, ...) a potřebného výpočetního výkonu. Tento

požadavek je dán tím, že velikost výsledné sítě (ať už co se týká počtu parametrů nebo hyperparametrů) hraje důležitou roli z pohledu požadavků na zdroje zařízení, na kterém poběží. V dnešní době mobilních a vestavěných zařízení s omezenými zdroji je tento požadavek zcela pochopitelný. Najít optimalizační metodu, která by byla schopna tento úkol zcela vyřešit, je ale prakticky nemožné, a tak je často nutné uchýlit se k různým heuristickým metodám. Ty jsou schopny na úkor přesnosti, úplnosti či efektivity nalézt takové řešení, které se co nejvíce blíží zadaným požadavkům. Jednu skupinu takovýchto metod tvoří evoluční výpočetní techniky [7], založené na poznacích ze světa biologické evoluce.

Cílem této práce je využití evolučních výpočetních technik jako optimalizační metody při automatizovaném návrhu nových architektur konvolučních neuronových sítí. Pro tyto účely byla zvolena speciální technika nazývaná kartézské genetické programování [8]. S využitím této techniky se lze zaměřit na hledání

21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40

41 zajímavých architektur z pohledu přesnosti, počtu pa- 91
42 rametrů či technických požadavků na výslednou síť. 92
43 Výsledkem této práce je implementace programu 93
44 pro automatizovaný návrh konvolučních neuronových 94
45 sítí, využívající uživatelem definované výpočetní pros- 95
46 tředky (v našem případě GPU), určených k řešení 96
47 specifických problémů.

48 2. Teoretické základy

49 2.1 Konvoluční neuronové sítě

50 Konvoluční neuronové sítě (Convolutional Neural Net- 97
51 works, zkráceně *CNN*) tvoří skupinu hlubokých neu- 98
52 ronových sítí, které se specializují na zpracování obra- 99
53 zových dat. Moderní konvoluční neuronové sítě se 100
54 skládají ze čtyř základních typů vrstev – *konvoluční*, 101
55 *aktivační*, *seskupující* a *plně propojené*. Jejich možné 102
56 umístění v architektuře CNN je uvedeno na obrázku 1. 103

57 Konvoluční vrstvy jsou zodpovědné za extrakci 104
58 užitečných rysů ze vstupních obrázků. K tomu využijí 105
59 operaci konvoluce, která realizuje aplikaci filtrů 106
60 (označovaných též jako jádro či kernel) na vstupní 107
61 obrázek. Ten má zpravidla několik kanálů (v případě 108
62 RGB obrázku 3), přičemž konvoluce je prováděna 109
63 zvlášť pro každý z nich.

64 Aktivační vrstvy plní důležitou roli, jelikož do 110
65 modelu umělých neuronových sítí vnáší prvek nelinea- 111
66 rity. Bez aktivačních vrstev by sebesložitější vícevrstvá 112
67 neuronová síť představovala jen o něco složitější, jed- 113
68 novrstvou neuronovou síť. Aktivační vrstvy jsou zpra- 114
69 vidla umístěny za lineární operace jako je konvoluce 115
70 či bázová funkce perceptronu. Nejpoužívanější ak- 116
71 tivační funkcí je usměrněná lineární jednotka (Recti- 117
72 fied Linear Unit, zkráceně *ReLU*) [9].

73 Účelem seskupujících vrstev je zachovat pouze 118
74 důležité informace z obrázku a zbavit se nepodstatných 119
75 detailů, jako je například umístění či natočení extraho- 120
76 vaných rysů. Seskupující vrstvy rovněž redukuje pros- 121
77 torové rozměry vnitřní reprezentace vstupního obrázku 122
78 (*downsampling*). Tato vrstva se nejčastěji umísťuje za 123
79 konvoluční, přičemž pracuje zvlášť pro každý vstupní 124
80 kanál.

81 Poslední část klasifikačních CNN tvoří plně propo- 125
82 jená neuronová vrstva, určená ke klasifikaci vstupních 126
83 obrázků do výstupních tříd. Před tuto část je zpravidla 127
84 nutné vložit zplošťující vrstvu, která transformuje vni- 128
85 trní tří dimenzionální reprezentaci obrázku (s dimen- 129
86 zemi výška, šířka a počet kanálů) do podoby vhodné 130
87 pro zpracování plně propojenou vrstvou.

88 2.2 Evoluční algoritmy

89 Evoluční algoritmy (zkráceně *EA*) jsou stochastické, 131
90 heuristicky založené algoritmy inspirované biologic-

91 kou evolucí. Jejich primárním úkolem je optimalizace 92
93 problémů, pro které neexistuje žádný efektivní algorit- 94
95 mus¹. Základní premisou konvergence EA je, že jedin- 96
97 ci, kteří splňují určité požadavky (jsou nadprůměrně 98
99 zdatní), budou mít potomky, a tak budou moci předat 100
101 své užitečné vlastnosti do další generace.

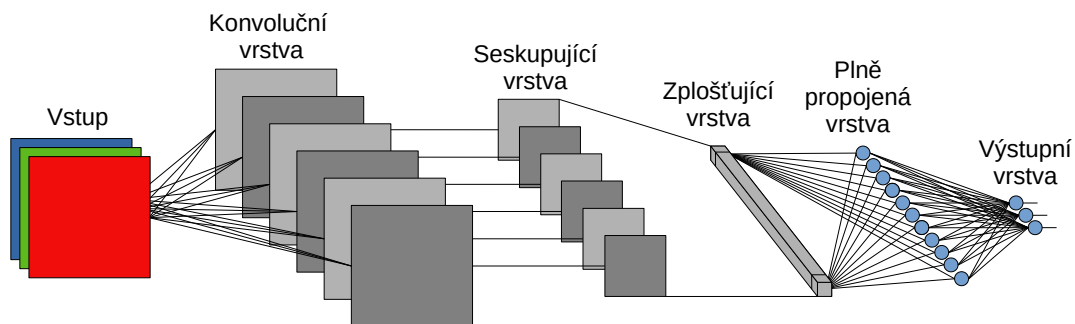
102 EA pracují nad množinou *jedinců*, kteří představují 103
104 nějaká validní řešení daného problému. Evoluce pro- 105
106 bíhá postupně po iteracích, nazývaných *generace*, do 107
108 té doby, dokud nejsou splněny *ukončující podmínky*. 109
110 Mezi ty patří například nalezení řešení s dostatečnou 111
112 kvalitou nebo vyčerpání maximálního počtu generací. 113
114 Množina jedinců, se kterými algoritmus pracuje v rám- 115
116 ci jedné generace, se nazývá *populace*. V každé ge- 117
118 neraci je s využitím operátoru *selektce* vybráno ně- 119
120 kolik jedinců, kteří se budou podílet na tvorbě nové 121
122 populace. Z vybraných jedinců jsou pomocí operátorů 123
124 *křížení* a *mutace* vyrobeni *potomci*, kteří tvoří popu- 124
125 lací následující generace. Jedinci jsou obvykle uloženi 125
126 v lineární struktuře nazývané *genotyp*. Jednotlivé části 126
127 genotypu se nazývají *geny*, přičemž konkrétní forma 127
128 genu se nazývá *alela*. Skutečná forma jedince se 128
129 nazývá *fenotyp*. 129
130

114 2.3 Kartézské genetické programování

115 Kartézské genetické programování (Cartesian Genetic 116
117 Programming, zkráceně *CGP*) [8] je speciální forma 118
119 genetického programování, ve které jednotliví jedinci 119
120 představují programy reprezentované acyklickými ori- 120
121 entovanými grafy (Directed Acyclic Graph, zkráceně 121
122 *DAG*). Uzly DAG jsou tvořeny výpočetními jednotkami, 122
123 uspořádanými ve dvoudimenzionální mřížce (odtud 123
124 název "*kartézské*"). 124

125 CGP má tři hlavní, uživatelem definované, para- 125
126 metry, počet řádků r , počet sloupců c a levels-back 126
127 l . První dva parametry udávají maximální počet vý- 127
128 početních uzlů, který je $r \times c$. Nastavením parametru 128
129 l uživatel udává, z kolika předchozích sloupců může 129
130 daný uzel brát své vstupy. Tento parametr tak ukládá 130
131 omezení při tvorbě propojení jednotlivých uzlů. Pokud 131
132 $l = 1$, tak může každý uzel brát vstup z výstupů uzlů 132
133 z levé strany sousedících sloupců nebo z primárních 132
134 vstupů. Změnou těchto parametrů lze docílit různých 133
135 topologií grafu. Speciálním případem je situace, kdy 133
136 $r = 1$ a $l = c$. Při tomto nastavení může vzniknout 134
137 libovolně propojený DAG, omezený pouze počtem 134
138 uzlů, který je roven c . 135
136
137
138

¹Efektivním algoritmem se z pohledu teorie složitosti rozumí algoritmus, který je schopný poskytnout řešení v polynomiálním čase.



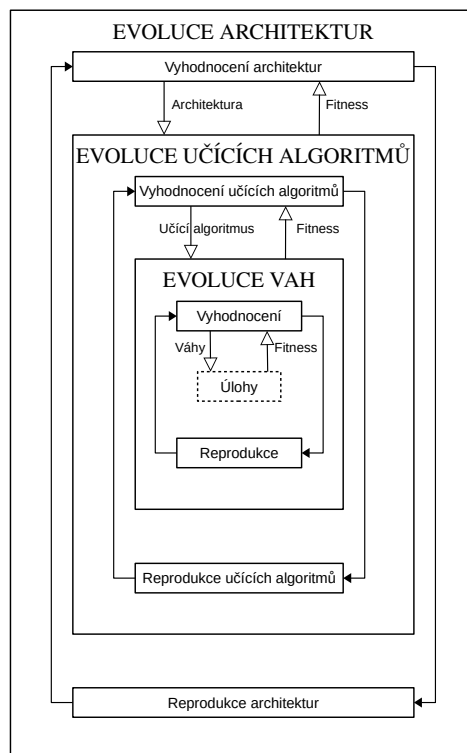
Obrázek 1. Ukázka jednoduché architektury CNN obsahující jednu konvoluční vrstvu, následovanou seskupující vrstvou, zplošťovací vrstvou, plně propojenou vrstvou a výstupní vrstvou.

139 Uživatel dále definuje počet vstupů a výstupů výpo-
 140 četní mřížky, která se skládá z funkčních uzlů, přičemž
 141 funkce uzlu je vybrána z uživatelem definované mno-
 142 žiny funkcí. Genotyp jedince má konstantní velikost
 143 a skládá se z genů jednotlivých uzlů a výstupních genů.
 144 Geny uzlů jsou dále rozděleny na funkční geny (udá-
 145 vající funkci uzlu) a propojovací geny, které stanovují,
 146 odkud bude uzel brát svůj vstup. Výstupní geny pak
 147 určují, z výstupů kterých uzlů bude tvořen výstup
 148 grafu.

149 Proces evoluce je u CGP řízen algoritmem in-
 150 spirovaným evoluční strategií ($\mu + \lambda$), nejčastěji ve
 151 formě $(1 + 4)$ [8]. V tomto jednoduchém algoritmu
 152 značí μ počet nejlepších jedinců, kteří jsou v každé
 153 generaci vybráni jako rodiče. Z nich je následně po-
 154 pomocí mutace vytvořeno λ potomků. V drtivé většině
 155 případů používá CGP pouze operaci mutace a nejčastěji
 156 pracuje s jedním rodičem a λ potomky v každé gene-
 157 racii.

158 2.4 Neuroevoluce

159 Neuroevoluce se zabývá využitím evolučních výpo-
 160 četních technik při návrhu umělých neuronových sítí.
 161 Evolučně navržené neuronové sítě se označují jako
 162 EANN (Evolutionary Artificial Neural Networks) [10].
 163 Proces evoluce je v kontextu tvorby umělých neuro-
 164 nových sítí využit v několika instancích, přičemž asi
 165 nejčastější je využití evoluce pro návrh architektury
 166 neuronové sítě. Zbylé dvě instance se zabývají evolucí
 167 vah nebo ještě obecněji návrhem učícího algoritmu.
 168 Evoluční algoritmy určené k evoluci struktur a vah
 169 neuronových sítí se označují jako TWEANN (Topo-
 170 logic and Weight Evolving Artificial Neural Network).
 171 Nejznámějším takovýmto algoritmem je genetický al-
 172 goritmus NEAT (NeuroEvolution of Augmenting Topo-
 173 logic), představený autorem Kennethem O. Stan-
 174 leyem v roce 2002 [11]. TWEANN algoritmy se nej-
 175 častěji dělí na konstruktivní a destruktivní [10]. Kon-
 176 struktivní algoritmy začínají s minimální funkční ar-
 177 chitekturou a v průběhu evoluce přidávají další prvky
 178 (uzly nebo spojení) do výsledné architektury. Des-



Obrázek 2. Univerzální schéma neuroevoluce. Převzato z [10].

179 truktivní algoritmy naopak začínají s архитектурou
 180 s maximálním množstvím prvků a v průběhu evoluce
 181 náhodné prvky eliminují. Obecné schéma neuroevoluce
 182 je uvedeno na obrázku 2.

3. Související práce

183 Tato sekce obsahuje krátké shrnutí a porovnání již e-
 184 xistujících prací, zabývajících se tematikou evolučního
 185 návrhu konvolučních neuronových sítí. Zvláštní po-
 186 zornost je zde věnována řešením zaměřujícím se na
 187 evoluční návrh struktur CNN s využitím CGP.
 188

189 Využití CGP při návrhu umělých neuronových sítí
 190 se označuje jako CGPANN (Cartesian Genetic Pro-
 191 gramming of Artificial Neural Networks) [12]. Jedná
 192 se o neuroevoluční algoritmus, který využívá přímého
 193 zakódování architektury, vah a funkcí do genotypu.

194 CGPANN tak představuje TWEANN algoritmus, který
 195 při procesu evoluce využívá výhradně operátoru mu-
 196 tace. Na rozdíl od genetických TWEANN algoritmů
 197 je CGPANN konstruktivní i destruktivní algoritmus.

198 Příklad úspěšného využití CGP při návrhu CNN
 199 je uveden v práci [13]. Její autoři vytvořili metodu
 200 nazvanou CGP-CNN, která používá CGP kódování
 201 jedinců, schopné reprezentovat architekturu a propo-
 202 jení výsledné CNN. Výhodnou tohoto kódování je
 203 jeho flexibilita, jelikož umožňuje reprezentaci různě
 204 hlubokých sítí i využití *skip* propojení. Mimo jed-
 205 noduché výpočetní uzly CGP mřížky, jako jsou kon-
 206 voluční a seskupující vrstva, autoři uvedli i složitější
 207 moduly, nazvané *ConvBlock* a *ResBlock*. Tyto mo-
 208 duly představují komplexnější výpočetní uzly vytvořené
 209 z jednodušších operací.

210 *ConvBlock* se skládá z konvoluce s krokem 1 a Re-
 211 LU aktivační funkce. Parametry tohoto modulu jsou
 212 velikost filtru (3×3 nebo 5×5) a počet výstupních
 213 kanálů (32, 64 nebo 128).

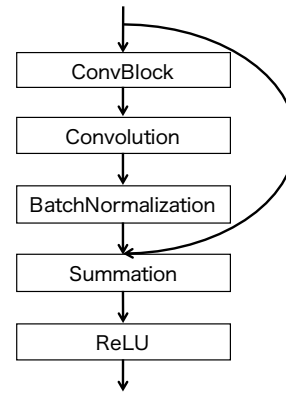
214 *ResBlock* je složen z konvoluce, konkatenace a Re-
 215 LU aktivační funkce. *Skip* propojení přeskakuje kon-
 216 voluční část tohoto modulu a přivádí nezměněný vstup
 217 modulu do konkatenčního uzlu, kde dochází ke konk-
 218 atenaci s výsledkem konvoluce. Na výsledek konkate-
 219 nace je následně použita aktivační funkce ReLU. Tento
 220 modul obsahuje parametry shodné jako u modulu *Con-*
 221 *vBlock*. Architektura modulu *ResBlock* je uvedena na
 222 obrázku 3.

223 Experimenty s touto metodou její autoři rozdělili
 224 na dva případy, podle toho, zda byly použity *Con-*
 225 *vBlock* nebo *ResBlock* moduly. Experimenty byly
 226 prováděny na datasetu CIFAR-10. Nejlepší vytvořená
 227 architektura v případě použití modulů *ConvBlock* do-
 228 sáhla přesnosti 93.25 % s 1.52M parametrů. Architek-
 229 tura vytvořená z modulů *ResBlock* dosáhla přesnos-
 230 ti až 94.02 % s 1.68M parametrů. Provedené ex-
 231 perimenty ukázaly, že metoda CGP-CNN je schopna
 232 vytvořit architektury porovnatelné s těmi ručně navr-
 233 ženými. V případě využití *ResBlock* modulů výsledná
 234 síť překonala i jedny z nejlepších ručně navržených
 235 sítí [13].

236 3.1 Srovnání existujících řešení

237 Tabulka 1 obsahuje srovnání přesnosti a počtu paramet-
 238 rů konvolučních neuronových sítí na datasetu CIFAR-
 239 10, navržených různými evolučními algoritmy. Pro
 240 srovnání obsahuje tabulka i současnou *state-of-the-*
 241 *art* odborně navrženou architekturu CNN nazvanou
 242 DenseNet [14].

243 Model NAS [4] využívá rekurentní neuronovou síť
 244 LSTM s posilovaným učením pro generování architek-
 245 tur CNN. Model nazvaný CNN-GA [15] je založen



Obrázek 3. Architektura modulu *ResBlock*. *BatchNormalization* provádí normalizaci hodnot po procesu konvoluce. Převzato z [13].

Srovnání algoritmů pro CIFAR-10			
Název modelu	Přesnost (%)	Počet parametrů	Doba běhu (GPU dny)
DenseNet [14]	96.54	25.6M	–
NAS [4]	93.99	2.5M	22 400
CGP-CNN [13]	94.02	1.68M	27
CNN-GA [15]	95.22	2.9M	35

Tabulka 1. Srovnání parametrů CNN, které byly získány pomocí různých přístupů (jeden GPU den znamená, že algoritmus strávil jeden den na jednom GPU).

na využití genetického programování pro návrh ar-
 246 chitektur CNN. Skvělých výsledků CNN-GA dosáhlo
 247 díky využití takzvaných *skip* propojení, které přeska-
 248 kují některé vrstvy v architektuře. Díky tomu dochází
 249 k lepší zpětné propagaci chyby do počátečních vrstev,
 250 a tak k rychlejší konvergenci. Posledním zkoumaným
 251 modelem v této sekci je CGP-CNN [13], využívající
 252 metodu CGP pro automatické hledání optimálních ar-
 253 chitektur CNN. 254

255 4. Navržené řešení

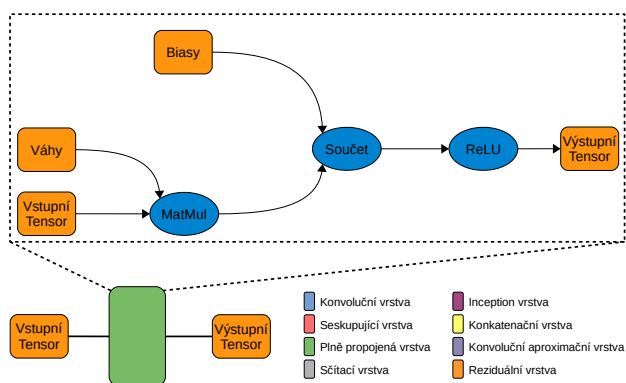
Pro automatizovaný návrh CNN architektur byla zv-
 256 olena metoda CGP, kde jednotlivé uzly CGP mřížky
 257 představují určité vrstvy CNN modelu. CGP tak hledá
 258 architekturu sítí CNN, přičemž učení vah je ponecháno
 259 na algoritmu stochastického gradientního sestupu s vy-
 260 užitím zpětného šíření chyby. 261

Pro práci s neuronovými sítěmi byla zvolena kni-
 262 hovna TensorFlow. Hlavními dvěma motivátory pro
 263 volbu této knihovny byla podpora akcelerace procesu
 264 učení na GPU a fakt, že výzkumná skupina *Evolvable*
 265 *Hardware* pracuje na rozšíření *tf-approximate*² kni-
 266 hovny TensorFlow o využití aproximačních jednotek. 267

Implementovaný program vytváří z výpočetních 268
 grafů knihovny TensorFlow jednotlivé vrstvy konvo- 269

²<https://github.com/ehw-fit/tf-approximate>

270 lučních neuronových sítí. Na obrázku 4 je uveden
 271 příklad implementace jedné vrstvy CNN, konkrétně
 272 plně propojené vrstvy, pomocí výpočetního grafu. Jak
 273 obrázek naznačuje, tak obdobným způsobem lze reali-
 274 zovat libovolnou vrstvu, jako je například konvoluční,
 275 seskupující, sčítací či konkatenáčn. Z těchto jednodu-
 276 chých vrstev lze vytvořit i složitější výpočetní celky,
 277 jako je reziduální blok (popsaný v sekci 3) nebo tzv.
 278 Inception modul [16].



Obrázek 4. Příklad implementace plně propojené vrstvy pomocí TensorFlow výpočetního grafu.

279 Funkční uzly CGP mřížky obsahují implementaci
 280 individuálních vrstev CNN architektury, jak je uvedeno
 281 na obrázku 5. Každý DAG tak představuje nějakou
 282 validní architekturu CNN. Nastavením počtu paramet-
 283 rů CGP mřížky může uživatel ovlivnit, jaké architek-
 284 tury bude navržený program vytvářet. Počet sloupců
 285 mřížky c udává maximální hloubku navrhovaných
 286 CNN. Pomocí počtu řádků r lze zase dát programu na
 287 výběr z několika alternativ funkčních uzlů v každém
 288 sloupci. Parametr l -back pak udává omezení navrho-
 289 vaných CNN z pohledu možných propojení jednotli-
 290 vých vrstev. Tento parametr rovněž umožňuje vytváře-
 291 ní užitečných skip propojení.

292 Dalším užitečným vstupem, který může uživatel
 293 programu poskytnout, je definice samotné výpočetní
 294 mřížky CGP. Vhodně rozmístěné funkční uzly ve vý-
 295 početní mřížce mohou totiž ulehčit hledání efektivních
 296 CNN architektur.

297 Primárním vstupem CGP je obrázek, vyjádřený
 298 jako 3D tensor, s dimenzemi výška, šířka a počet
 299 kanálů vstupního obrázku. CGP pracuje s tensor-
 300 y s hodnotami typu float. Primárním výstupem CGP je
 301 1D tensor (vektor), vyjadřující příslušnost vstupního
 302 obrázku do jedné z výstupních tříd. Výstupní tensor
 303 pro tyto účely používá kódování 1 z N .

304 4.1 Mutace

305 Implementovaný program pro zakódování jedinců ne-
 306 využívat klasické kódování chromozomu jako sekven-
 307 ce celých čísel, jak je tomu u většiny genetických al-

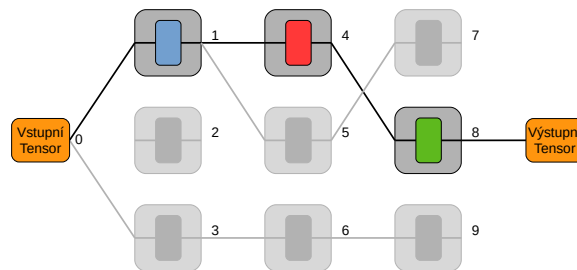
goritmů. V tomto případě je jedinec reprezentován 308
 přímo jako DAG, tedy dvojice $G = (U, E)$, kde U 309
 značí množinu uzlů a E množinu hran. Mutace je 310
 prováděna na grafu G , přičemž při mutaci dochází 311
 pouze ke změně propojení uzlů. V terminologii ge- 312
 netického programování to znamená, že mutace mění 313
 pouze propojovací a výstupní geny chromozomu. Fun- 314
 kční geny zůstávají nezměněny. 315

Mutace probíhá tak, že je náhodně vybrán uzel 316
 k mutaci $n \in U$. Pro uzel n jsou odstraněny všechny 317
 vstupní hrany, tedy $E = E \setminus \{(x, n) \mid x \in U\}$. Následně 318
 je spočítána množina 319

$$L(n) = \{u \mid u \in U \text{ a } u \text{ je uzel respektující parametr } l\text{-back pro vybraný uzel } n\}$$

Z množiny $L(n)$ je následně náhodně vybrán jeden 320
 uzel $m \in L(n)$. Do množiny hran E je poté přidána 321
 nová hrana (m, n) , tedy $E = E \cup (m, n)$. Pokud uzel 322
 n reprezentuje funkci s vyšší aritou než je 1, tak je 323
 proces přidání nové hrany opakován. 324

Mutace jedince je prováděna do té doby, dokud 325
 není vytvořen aktivní podgraf spojující vstupní uzel 326
 s výstupním. 327



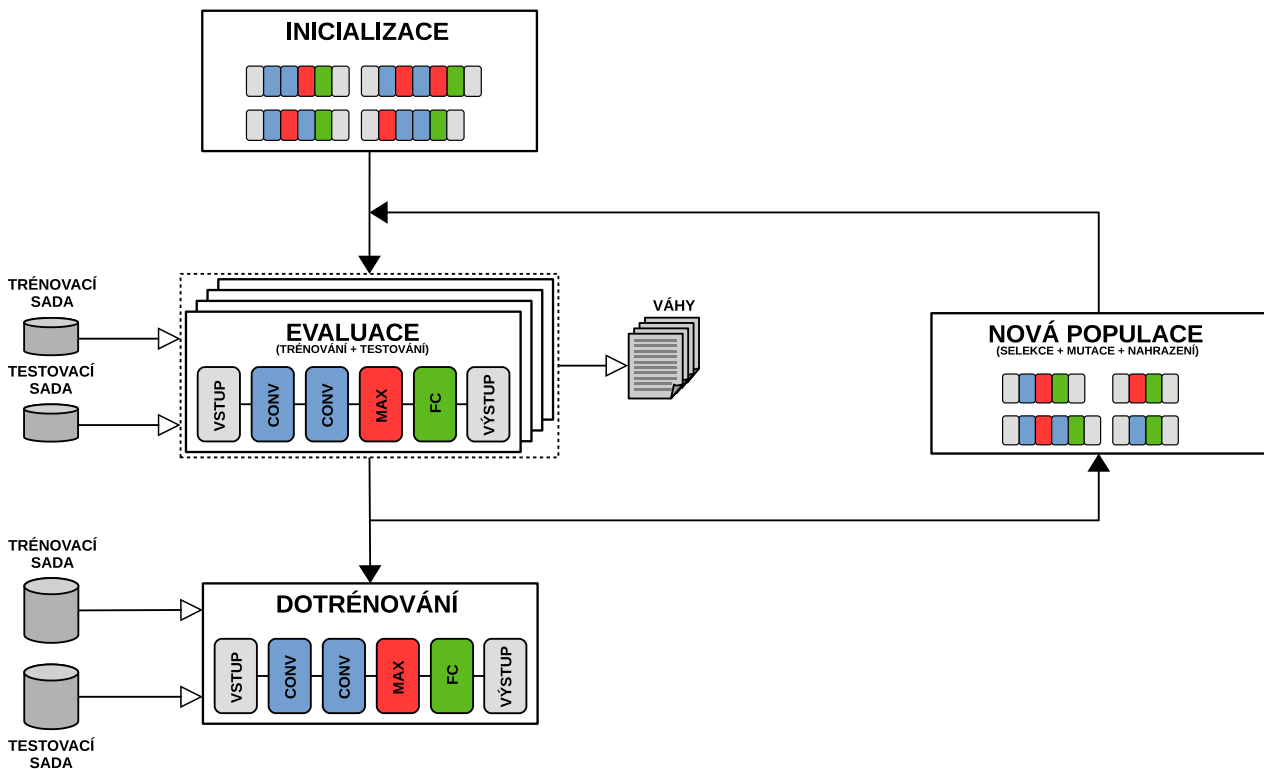
Obrázek 5. Ukázka DAG v CGP mřížce představující nalezené řešení, skládající se z konvoluční vrstvy (modrá), seskupující vrstvy (červená) a plně propojené vrstvy (zelená).

4.2 Prohledávací algoritmus 328

Implementovaný program využívá prohledávací algo- 329
 ritmus založený na evoluční strategii (1 + 4). Diagram 330
 prohledávacího algoritmu je uveden na obrázku 6. 331

Na úvod jsou vygenerovány 4 kandidátní řešení. 332
 Algoritmus dále vstupuje do smyčky, která provádí 333
 evaluaci každého jedince. Ta se skládá z trénování jed- 334
 inců pomocí náhodně vybrané podmnožiny trénovací 335
 datové sady. Následně testování jedinců je opět pro- 336
 váděno na náhodně vybrané podmnožině testovací 337
 datové sady. Všichni jedinci v určité generaci jsou 338
 trénováni a testováni na stejných, náhodně vybraných 339
 podmnožinách originálních datových sad. Tento postup 340
 byl převzat z [13]. 341

Během evaluace je jedincům přiřazena fitness hod- 342
 nota podle jejich přesnosti a počtu parametrů. Výpočet 343



Obrázek 6. Diagram prohledávacího algoritmu. Jak je zde znázorněno, tak všichni jedinci v každé generaci jsou trénováni a testováni s využitím náhodně vybrané podmnožiny originální datové sady. Naučené váhy každého jedince jsou uloženy do souboru. Po skončení algoritmu je nejlepší jedinec (do)trénován na originální datové sadě.

344 fitness je proveden pomocí vzorce

$$f = a \times \left(\frac{k}{\log(p) + 1} + 1 \right), \quad (1)$$

345 kde a značí přesnost, p počet parametrů a koeficient
346 k udává vliv počtu parametrů na výslednou fitness
347 hodnotu. Tato fitness funkce byla převzata z [17].

348 Z ohodnocených jedinců je poté vybrán nejlepší
349 jedinec – rodič. Pokud je v současné populaci více jed-
350 inců s nejlepší fitness hodnotou, tak je vybrán jedinec,
351 který ještě nebyl rodičem. V případě, že je takovýchto
352 jedinců více, je zvolen náhodný z nich.

353 S využitím operátoru mutace, popsaném v pod-
354 sekci 4.1, jsou z rodiče vytvořeni potomci, kteří tvoří
355 populaci následující generace.

356 Hlavní smyčka prohledávacího algoritmu končí,
357 až je dosaženo maximálního počtu generací. Nejlepší
358 nalezené řešení je následně (do)trénováno s využitím
359 originální datové sady.

360 4.3 Uložení nejlepších vah

361 Z popisu prohledávacího algoritmu plyne, že v každé
362 generaci jsou všichni jedinci podrobni trénování před
363 tím, než je určena jejich přesnost. Do další generace
364 se však dostanou jen někteří jedinci, z čehož plyne,
365 že trénování některých jedinců bylo zbytečné a jejich

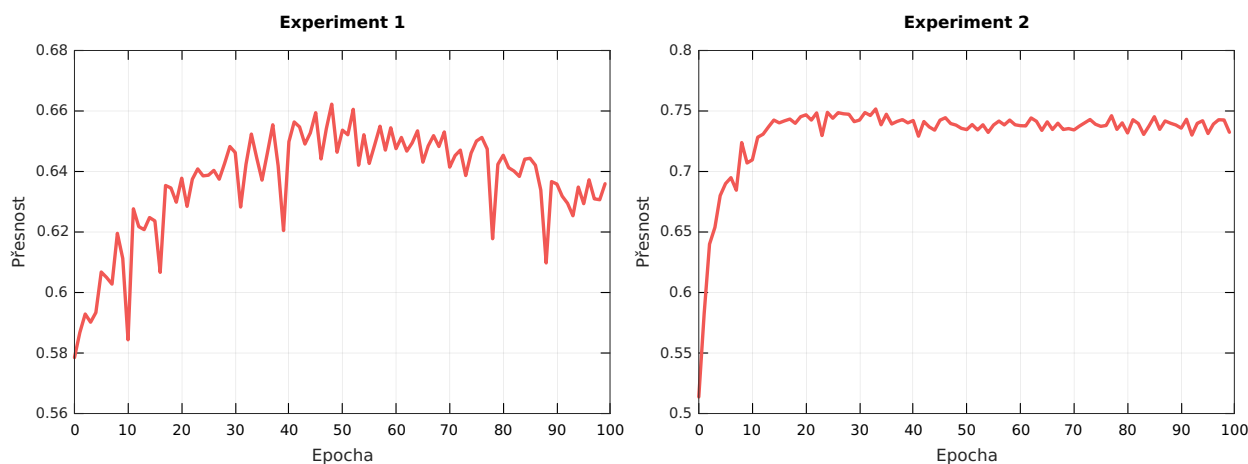
naučené váhy se zahodí. Další negativní vlastností to- 366
hoto přístupu je to, že noví jedinci budou muset začínat 367
od začátku, s náhodně vygenerovanými váhami. 368

Implementovaný program tento problém řeší tak, 369
že každý funkční uzel v CGP výpočetní mřížce si pa- 370
matuje naučené váhy nejlepšího jedince, který tento 371
uzel využíval. Noví jedinci si pak při trénování načtou 372
tyto váhy a nebudou tak muset začínat od znovu. Je- 373
likož jsou ale rozměry váhového tensoru každého funk- 374
kčního uzlu určeny tím, ke kterému funkčnímu uzlu je 375
připojen, tak se mohou rozměry váhového tensoru pro 376
ten samý uzel lišit. V tomto případě je váhový tensor 377
buď oříznut nebo doplněn náhodnými hodnotami tak, 378
aby odpovídal potřebným rozměrům. 379

5. Experimenty

380

Experimenty byly prováděny na datové sadě CIFAR- 381
10, která se skládá z 60 000 barevných obrázků o ve- 382
likosti 32×32 , rozdělených do 10 tříd. Datová sada 383
je dále rozdělena na trénovací sadu obsahující 50 000 384
vzorků a testovací sadu o velikosti 10 000 vzorků. Ex- 385
perimenty byly spuštěny na stroji se čtyřmi grafickými 386
kartami NVIDIA GTX 1080 (Pascal), 8GB RAM, 387
které byly použity pro akceleraci procesu učení CNN. 388
Všechny ostatní výpočty byly prováděny na CPU. 389



(a) Přesnost během finálního dotrénování nejlepšího řešení, nalezeného v experimentu 1.

(b) Přesnost během finálního dotrénování nejlepšího řešení, nalezeného v experimentu 2.

Obrázek 7. Grafy experimentů.

390 5.1 Experiment 1

391 První experiment měl za cíl navrhnout architekturu
392 CNN s použitím jednoduchých výpočetních bloků,
393 jako jsou konvoluční, seskupující, konkatenáční, sčítací
394 a plně propojené vrstvy.

395 5.1.1 Nastavení programu

396 CGP mřížka obsahovala 6 řádků a 31 sloupců s parametrem l-back nastaveným na 8. Mřížka byla vytvořena tak, že sloupce blíže primárnímu vstupu obsahovaly konvoluční vrstvy s větší velikostí filtru a menším počtem výstupních kanálů. Čím blíže se pak sloupce blížily výstupu CGP mřížky, tím se velikosti filtrů zmenšovaly a počty výstupních kanálů zvyšovaly. Tento postup byl zvolen z toho důvodu, že hlubší vrstvy si mohou dovolit mít více výstupních kanálů, jelikož pracují s menšími vstupy. Konvoluční sloupce CGP mřížky byly příležitostně proloženy sloupci se seskupujícími, konkatenáčními a sčítacími vrstvami. Poslední sloupec se skládal z plně propojených vrstev s počtem neuronů 64, 128, 256, 512, 1024 a 2048.

410 Počet generací prohledávacího algoritmu byl nastaven na 50, přičemž algoritmus v každé generaci pracoval s populací velikosti 4. V každé generaci bylo z trénovací datové sady náhodně vybráno 10 000 vzorků a z testovací datové sady 5 000 náhodných vzorků. Pomocí této mini-trénovací sady bylo provedeno trénování každého jedince po dobu 10 epoch. Přesnost každého kandidátního řešení pak bylo určeno evaluací na mini-testovací sadě. Parametr k fitness funkce 1 byl nastaven na 0.5.

420 Po skončení prohledávacího algoritmu bylo nejlepší nalezené řešení (do)trénováno na plné trénovací sadě po dobu 100 epoch. Výsledná přesnost pak byla spočítána pro originální testovací sadu.

5.2 Experiment 2

Cílem druhého experimentu bylo navrhnout architekturu CNN s využitím ResBloku, popsaném na obrázku 3.

5.2.1 Nastavení programu

CGP mřížka obsahovala 6 řádků a 15 sloupců s parametrem l-back nastaveným na 4. Jednotlivé sloupce CGP mřížky obsahovaly ResBloky, konvoluční, seskupující, sčítací a konkatenáční vrstvy s různými parametry. Poslední sloupec se skládal z plně propojených vrstev s počtem neuronů 64, 128, 256, 512, 1024 a 2048.

Nastavení prohledávacího algoritmu tohoto experimentu bylo shodné s nastavením v experimentu 1.

5.3 Výsledky

Přesnosti nejlepších nalezených řešení obou experimentů jsou ukázány na obrázku 7.

Výsledkem prvního experimentu bylo řešení s přesností 64.5 % a počtem parametrů 146178. Graf 7a zobrazuje průběh učení nejlepšího řešení při finálním dotrénování, které trvalo 100 epoch. Jak je vidět, tak přesnost dosáhla až k 66 % a přibližně od 50. epochy se snižovala. To ukazuje, že při procesu učení začalo docházet k přeučení (overfitting), kdy model začal ztrácet schopnost generalizace.

Výsledkem druhého experimentu bylo řešení s přesností 74.5 % a počtem parametrů 475722. Graf 7b zobrazuje průběh učení nejlepšího řešení při finálním dotrénování, které trvalo 100 epoch. Jak je vidět, tak přesnost dosáhla až k 75 % a přibližně od 30. epochy se opět začala pozvolna snižovat. To opět ukazuje na přeučení.

Cílem této práce bylo navrhnout a implementovat neuroevoluční algoritmus pro automatizovaný návrh architektur konvolučních neuronových sítí. Toho bylo dosaženo použitím speciální techniky CGP. Pro práci s CNN byla zvolena knihovna TensorFlow, umožňující akceleraci výpočtů na GPU. Pro ověření funkčnosti implementovaného programu byly provedeny dva experimenty na datové sadě CIFAR-10, které dosáhly přesnosti 64.5 % s 146178 parametry a 74.5 % s počtem parametrů 475722.

Výsledkem této práce je program, implementující evoluční návrh architektur konvolučních neuronových sítí. Implementovaný program rovněž umožňuje uživateli specifikovat si různé požadavky a omezení na navrhované CNN. Uživatel tak může programu říci, jaké výpočetní jednotky může používat nebo jaká může být minimální či maximální hloubka sítě. Z praktického pohledu je tato možnost velmi užitečná, uvážíme-li, že bychom například hledali architekturu CNN pro nějaké zařízení s omezenými výpočetními prostředky.

6.1 Pokračování práce

Styl, jakým je program navržen a implementován podporuje možnost rozšíření CGP mřížky o nové, výkonnější funkční bloky. Implementovaný program může být rovněž použit pro návrh takovýchto bloků. Další alternativou je inicializace CGP nejlepšími známými CNN a hledání kompromisu mezi přesností a velikostí sítě.

Jiným podstatným rozšířením tohoto programu je využití aproximačních jednotek, které na úkor přesnosti snižují požadavky na příkon nebo plochu na čipu. Rozšíření *tf-approximate* implementuje tyto aproximační jednotky, ale nepodporuje proces učení. Tyto jednotky by ale mohly být lehce integrovány do implementovaného programu tak, že v procesu evaluace jedině by trénování probíhalo na přesných jednotkách, zatímco testování na těch aproximovaných. Výsledkem programu by pak byla naučená architektura CNN, která by mohla být na čipu implementována aproximačními jednotkami a tak ušetřit cenné zdroje vestavěného zařízení.

Poděkování

Rád bych poděkoval svému vedoucímu práce panu profesoru Ing. Lukáši Sekaninovi, Ph.D. za odborné vedení, konzultace, čas a cenné rady při tvorbě této práce. Dále bych rád poděkoval doktoru Ing. Vojtěchu Mrázkovi a docentu Ing. Jiřímu Jarošovi, Ph.D. za zprostředkování přístupu na školní výpočetní server.

- [1] Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. *Neural Architecture Search: A Survey*. *Journal of Machine Learning Research*, 20(55):1–21, 2019.
- [2] Md Ashiqur Rahman. *Neural Architecture Search (NAS) - The Future of Deep Learning*, June 2019. URL: <https://towardsdatascience.com/neural-architecture-search-nas-the-future-of-deep-learning-c99356351136>.
- [3] Kenneth O. Stanley, Jeff Clune, Joel Lehman, and Risto Miikkulainen. *Designing neural networks through neuroevolution*. *Nature Machine Intelligence*, 1(1):24–35, 2019.
- [4] Barret Zoph and Quoc V. Le. Neural architecture search with reinforcement learning. *International Conference on Learning Representations*, 2017.
- [5] Lewei Yao, Hang Xu, Wei Zhang, Xiaodan Liang, and Zhenguo Li. *SM-NAS: Structural-to-Modular Neural Architecture Search for Object Detection*, 2019. URL: <https://arxiv.org/abs/1911.09929>.
- [6] Ning Zhu. *Neural Architecture Search for Deep Face Recognition*, 2019. URL: <http://arxiv.org/abs/1904.09523>.
- [7] Patrick Siarry, Alain Petrowski, and Sana Ben Hamida. *Metaheuristics*, chapter Evolutionary Algorithms. Springer International Publishing, 2016.
- [8] J. F. Miller. *Cartesian Genetic Programming*, pages 17–34. Natural Computing Series. Springer Berlin Heidelberg, 2011.
- [9] Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on International Conference on Machine Learning, ICML'10*, page 807–814, Madison, WI, USA, 2010. Omnipress.
- [10] Xin Yao. *Evolving artificial neural networks*. *Proceedings of the IEEE*, 87(9):1423–1447, September 1999.
- [11] Kenneth O. Stanley and Risto Miikkulainen. Evolving neural networks through augmenting topologies. *Evolutionary Computation*, 10(2):99–127, 2002.
- [12] Maryam Mahsal Khan and Gul Muhammad Khan. A novel neuroevolutionary algorithm:

- 553 Cartesian genetic programming evolved arti-
554 ficial neural network (cgpann). In *Proceedings of*
555 *the 8th International Conference on Frontiers of*
556 *Information Technology*, New York, NY, USA,
557 2010. Association for Computing Machinery.
- 558 [13] Masanori Suganuma, Shinichi Shirakawa, and
559 Tomoharu Nagao. A genetic programming ap-
560 proach to designing convolutional neural network
561 architectures. In *Proceedings of the Genetic and*
562 *Evolutionary Computation Conference, GECCO*
563 *'17*, page 497–504, New York, NY, USA, 2017.
564 Association for Computing Machinery.
- 565 [14] G. Huang, Z. Liu, L. v. d. Maaten, and K. Q.
566 Weinberger. Densely connected convolutional
567 networks. In *2017 IEEE Conference on Com-*
568 *puter Vision and Pattern Recognition (CVPR)*,
569 pages 2261–2269, July 2017.
- 570 [15] Yanan Sun, Bing Xue, Mengjie Zhang, and
571 Gary G. Yen. Automatically designing CNN
572 architectures using genetic algorithm for image
573 classification. 2018.
- 574 [16] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre
575 Sermanet, Scott Reed, Dragomir Anguelov, Du-
576 mitru Erhan, Vincent Vanhoucke, and Andrew
577 Rabinovich. Going deeper with convolutions.
578 *2015 IEEE Conference on Computer Vision and*
579 *Pattern Recognition (CVPR)*, pages 1–9, 2015.
- 580 [17] Filip Badáň and Lukáš Sekanina. Optimizing
581 convolutional neural networks for embedded sys-
582 tems by means of neuroevolution. In *Theory*
583 *and Practice of Natural Computing*, pages 109–
584 121. Springer International Publishing, Listopad
585 2019.