

Uživatelská dokumentace

Tato dokumentace obsahuje návod jak spustit ukázkou nové implementace indexu HNSW, která je součástí bakalářské práce na téma "Aproximace KNN problému".

Všechny cesty uvedené v tomto souboru jsou relativní k cestě složky, která obsahuje tuto dokumentaci.

Potřebné programy

- Docker
- Překladač C++17
- Python 3.9

Ukázka

1. Ujistěte se, že služba Docker je zapnutá.
2. Spustíte skript `RUNME.py` pomocí interpretu Python verze 3.9. Na Windows například takto:

```
1 py -3.9 RUNME.py
```

Tento skript provede následující.

- Vytvoří virtuální prostředí interpretu Python ve složce `.venv`.
- Vytvoří nativní C++ řešení ve složce `src/cmakeBuild`.
- Spustí porovnání původní a nové implementace indexu HNSW nad malými kolekcemi se 100 000 prvky.
- Otevře stránku s výsledky v jedné kartě internetového prohlížeče.
- Otevře stránku s dokumentací indexu ve druhé kartě.

Výsledky srovnání

Výsledky jsou zaznačeny do grafů, které zobrazuje webová stránka. Ta je generována ve složce `src/website`. Původní implementace je v grafech označována slovem `original`, nová implementace slovem `new`.

Virtuální prostředí

Pokud není uvedeno jinak, skripty uvnitř složky `src/scripts` vždy spouštějte pomocí vygenerovaného virtuálního prostředí. Prostředí aktivujete pomocí aktivačního skriptu ve složce `.venv/Scripts`. Výběr skriptu závisí na použitém OS a interpretu.

OS	INTERPRET	CESTA K AKTIVAČNÍMU SKRIPTU
Linux		<code>./venv/Scripts/activate</code>
Windows	Batch	<code>.\venv\Scripts\activate.bat</code>
Windows	Powershell	<code>.\venv\Scripts\Activate.ps1</code>

Seznam skriptů

Následuje seznam skriptů ve složce `src/scripts`, které umožňují uživateli provést více operací než úvodní skript.

NÁZEV SKRIPTU	STRUČNÝ POPIS SKRIPTU
<code>buildProject</code>	Vytvoří virtuální prostředí, nativní C++ řešení a jeho Python rozhraní.
<code>clean</code>	Odstraní vygenerované soubory a vrátí projekt do původního stavu.
<code>datasetGenerator</code>	Vygeneruje datové soubory pro debugování.
<code>datasetToText</code>	Převede datový soubor do textového formátu.
<code>formatCMakeTemplates</code>	Vygeneruje CMakeLists.txt.
<code>generateTables</code>	Vygeneruje LaTeX tabulky podobné těm, které jsou v bakalářské práci.
<code>latexTable</code>	Vygeneruje LaTeX tabulku na základě výsledků srovnání.
<code>runBenchmarks</code>	Spustí srovnání, vygeneruje a otevře webovou stránku s výsledky.
<code>runRecallTable</code>	Postaví nový index a zobrazí tabulku závislosti přesnosti na parametru vyhledávání <code>ef_{search}</code> .
<code>SIMDCapability</code>	Zobrazí SIMD rozšíření instrukční sady procesoru, která jsou k dispozici.

Podrobný popis skriptů

U každého skriptu je uveden jeho účel, parametry a příklad spuštění. Pokud skript obsahuje alespoň jeden parametr, pak použitím parametru `--help` nebo `-h` zobrazíte nápovědu v anglickém jazyce.

buildProject

Tento skript lze spustit bez virtuálního prostředí.

Vytvoří virtuální prostředí interpretu Python, stáhne potřebné softwarové balíčky, vygeneruje nativní C++ řešení pro knihovnu nového indexu, vytvoří rozhraní v jazyce Python pro nový index a otestuje funkčnost tohoto indexu spuštěním skriptu `runRecallTable`.

PARAMETR, ZKRATKA	VÝZNAM
<code>--clean, -c</code>	Vrátí projekt do původního stavu před jeho opětovným sestavením.
<code>--cleanResults, -r</code>	Pokud je <code>--clean</code> nastaven, odstraní naměřené výsledky.
<code>--ignorePythonVersion, -i</code>	Umožňuje spustit skript s libovolnou verzí interpretu Python. Skript poté nemusí fungovat správně.

Příklad spuštění:

```
1 py -3.9 buildProject.py --clean --cleanResults
```

clean

Tento skript lze spustit bez virtuálního prostředí.

Odstraní datové soubory pro debugování, C++ nativní řešení a Python rozhraní. Pokud je spuštěn mimo virtuální prostředí, pak odstraní toto prostředí. Naměřené výsledky odstraněny nebudou, pokud o to uživatel nepožádá.

PARAMETR, ZKRATKA	VÝZNAM
<code>--results, -r</code>	Odstraní naměřené výsledky srovnání, vygenerované grafy a tabulky.

Příklad spuštění:

```
1 py clean.py --results
```

datasetGenerator

Vygeneruje datové soubory pro debugování uvedené v konfiguračním souboru `src/config/debugDatasets.json`. O konfiguraci tohoto skriptu se více dočtete v kapitole `Datové soubory` níže v této dokumentaci.

Příklad spuštění:

```
1 py datasetGenerator.py
```

datasetToText

Převěde vybraný datový soubor ze složky `src/data` do textového formátu.

PARAMETR, ZKRATKA	VÝZNAM
--name, -n	Název datového souboru bez přípony. Pokud není uveden, výchozím souborem je <code>angular-small</code> .

Příklad spuštění:

```
1 py datasetToText --name euclidean-medium
```

formatCMakeTemplates

Vygeneruje soubor `src/index/CMakeLists.txt` a doplní do něj správnou definici maker tak, aby došlo pouze ke kompilaci těch funkcí, pro které je k dispozici vhodné SIMD rozšíření instrukční sady procesoru.

Příklad spuštění:

```
1 py formatCMakeTemplates.py
```

generateTables

Vygeneruje LaTeX tabulky podobné těm, které jsou v bakalářské práci, ale pouze v případě, že jsou pro ně dostupné naměřené výsledky. Tyto výsledky lze získat spuštěním následujících příkazů. Avšak tato měření mohou trvat více než 12 hodin.

```
1 py runBenchmarks.py -a ..\config\heuristic.yaml
  ..\config\naive.yaml -d lastfm-64-dot -r 5
2 py runBenchmarks.py -a ..\config\heuristic.yaml
  ..\config\prefetch.yaml -d glove-50-angular -r 5
3 py runBenchmarks.py -a ..\config\original.yaml
  ..\config\prefetch.yaml -d sift-128-euclidean -r 5
4 py generateTables.py
```

Vygenerované tabulky jsou dostupné ve složce `src/figures`.

latexTable

Vygeneruje jednu LaTeX tabulku na základě výsledků srovnání implementací.

PARAMETR, ZKRATKA	VÝZNAM
--algorithms, -a	Vyžadován. Seznam implementací oddělený mezerami.
--dataset, -d	Vyžadován. Název datového souboru.
--label, -la	Identifikátor tabulky.
--legend, -le	Názvy implementací v tabulce. Pokud není uveden, budou použity původní názvy.
--output, -o	Vyžadován. Cesta k výstupnímu souboru.
--percent, -p	Přidá do tabulky sloupec s procentuálním rozdílem časů stavby.
--recompute, -r	Znovu vypočítá výkonnostní metriky z naměřených výsledků. Tato operace může trvat více než 10 minut.

Příklad spuštění:

```
1 py latexTable.py -a new-prefetch original -d sift-128-euclidean -
  le "Nová impl." "Původní impl." -o ..\figures\table.tex -p
```

runBenchmarks

Před spuštěním se ujistěte, že je služba Docker zapnutá.

Spustí srovnání implementací v jednom nebo více Docker kontejnerech, vypočítá výkonnostní metrika, vygeneruje webovou stránku s výsledky a otevře ji v nové kartě internetového prohlížeče. Kód vygenerované stránky lze poté najít ve složce `src/website` a můžete ji opětovně zobrazit otevřením souboru `index.html`.

PARAMETR, ZKRATKA	VÝZNAM
<code>--algoDefPaths, -a</code>	Vyžadován. Seznam cest ke konfiguračním souborům oddělených mezerami. O konfiguraci se více dočtete v kapitole <code>Konfigurace srovnání</code> .
<code>--datasets, -d</code>	Vyžadován*. Seznam datových souborů oddělených mezerami.
<code>--datasetsPath, -p</code>	Vyžadován*. Cesta k textovému souboru se seznamem datových souborů.
<code>--force, -f</code>	Spustí již provedená měření znovu.-
<code>--runs, -r</code>	Počet opakování měření. Výchozí hodnota je 1.
<code>--workers, -w</code>	Počet paralelně spuštěných Docker kontejnerů. Výchozí hodnota je 1.

Datové soubory využité ke srovnání nejsou ty samé, které jsou využívány k debugování. Jejich seznam najdete v kapitole `Testované datové soubory`.

* Pouze jeden z parametrů označených hvězdičkou by měl být uveden.

Příklad spuštění:

```
1 py runBenchmarks.py -a ..\config\noBit.yaml -f -p  
   ..\config\datasets.txt -r 5 -w 2
```

runRecallTable

Postaví index nové implementace a vyhledá v něm nejbližší sousedy s různými hodnotami parametru vyhledávání `efsearch`. Poté vypíše tabulku závislosti přesnosti vyhledávání na tomto parametru. Konfigurace tohoto skriptu se nachází v souboru `src/config/recallTable.json` a více se o ní dočtete v kapitole `Konfigurace recallTable`.

Příklad spuštění:

```
1 py runRecallTable.py
```

SIMDCapability

Zobrazí SIMD rozšíření instrukční sady procesoru, která jsou k dispozici. Využívána ostatními skripty pro vygenerování správných maker v jazyce C++.

Příklad spuštění:

```
1 py SIMDCapability.py
```