

Zarovňávanie nanopórových čítaní ku grafom

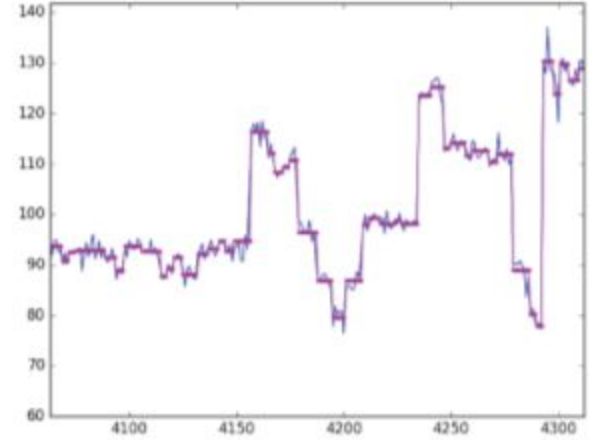
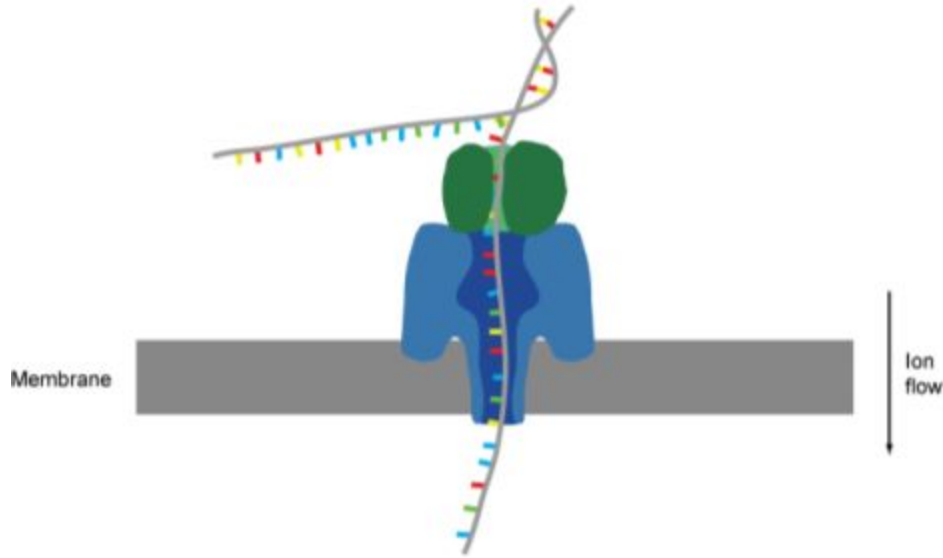
Názov (EN): Alignment of Nanopore Reads to Graphs

Názov (SK): Zarovňávanie nanopórových čítaní ku grafom

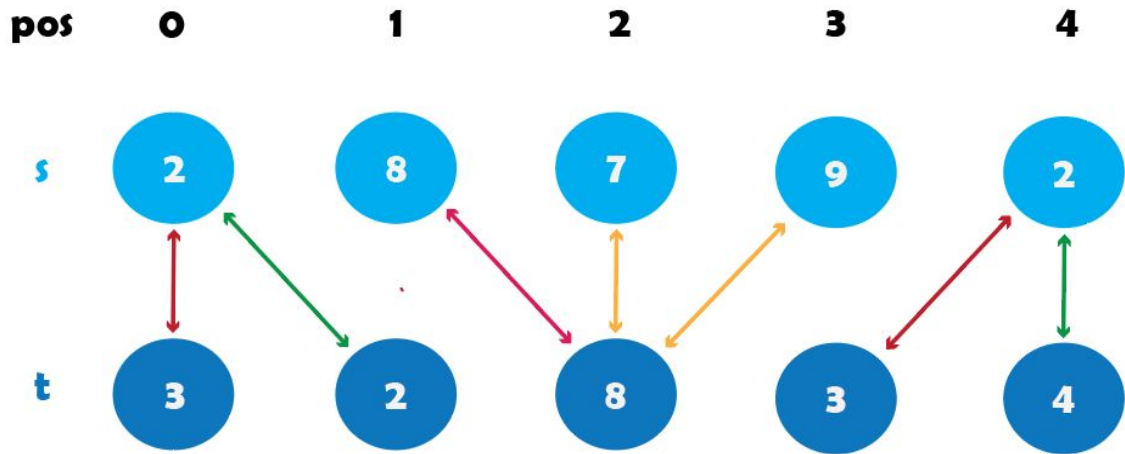
Študent: Roman Beňo

Školiteľ: doc. Mgr. Tomáš Vinař, PhD.

I. Stručný popis problematiky



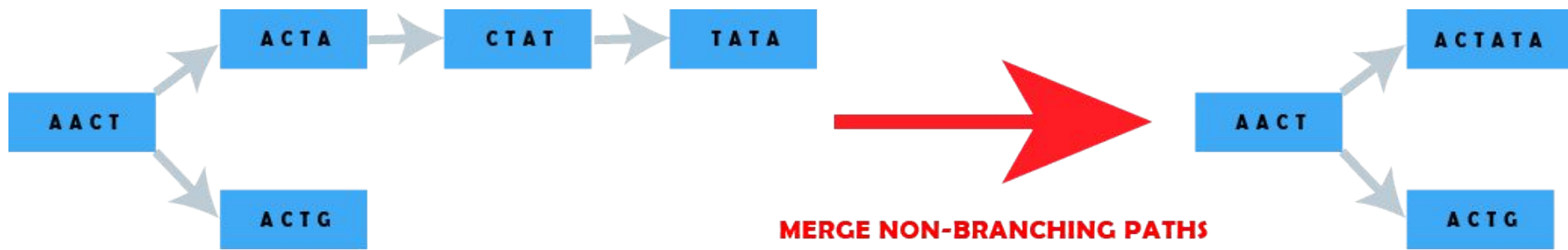
Nanopórové sekvenovanie, MinION, squiggle (“čmáranec”), špecifiká výstupných dát, base-calling



pos	-	1	2	3	4	5
-	val	3	2	8	3	4
1	2	1	1	7	8	10
2	8	6	7	1	6	10
3	7	10	11	2	5	8
4	9	16	17	3	8	10
5	2	17	16	9	4	6

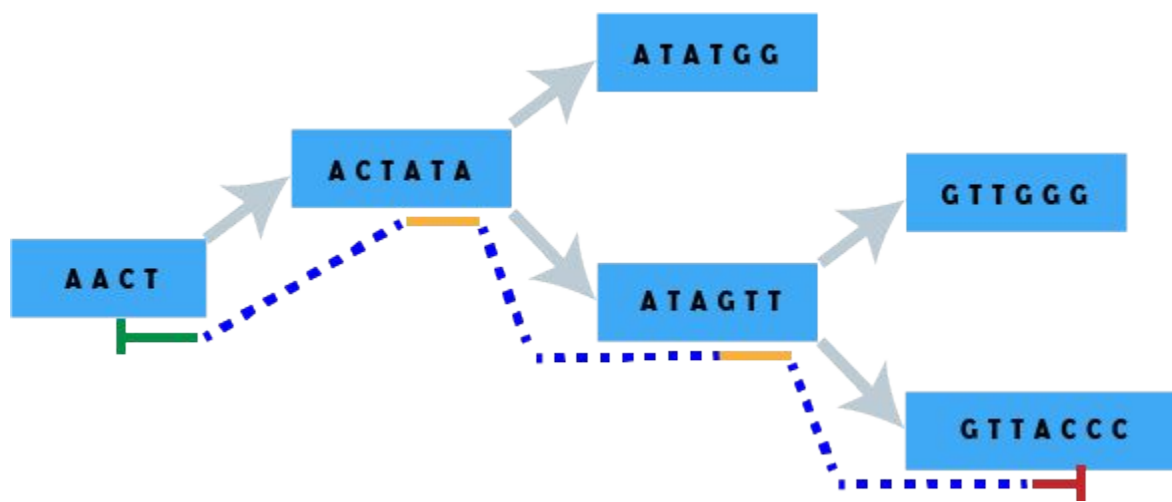
The table shows a cost matrix for DTW. A yellow circle with the number 0 is at the top-left corner (row 1, column 1). Red arrows point from this circle to the cells (1,2), (2,1), (2,2), (3,1), (3,2), (4,1), (4,2), (5,1), and (5,2). Green arrows point from the cells (1,2) to (1,3), (1,3) to (1,4), (1,4) to (1,5), (2,1) to (2,2), (2,2) to (2,3), (3,2) to (3,3), (3,3) to (3,4), (4,3) to (4,4), (4,4) to (4,5), and (5,4) to (5,5).

Dynamic time warping - penalizácia/cena (penalty), funkcia na výpočet ceny (cost function), minimum warping path

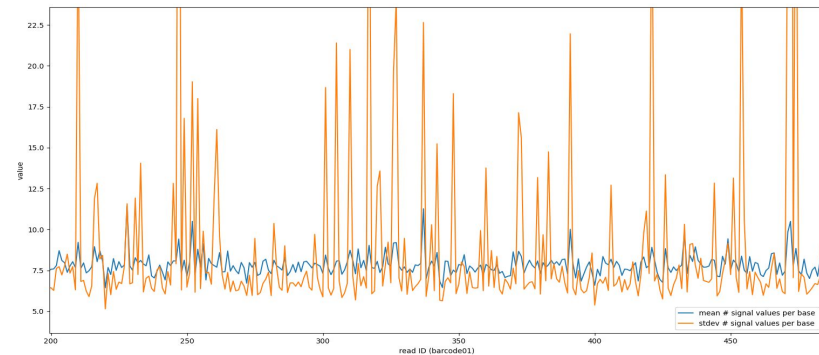
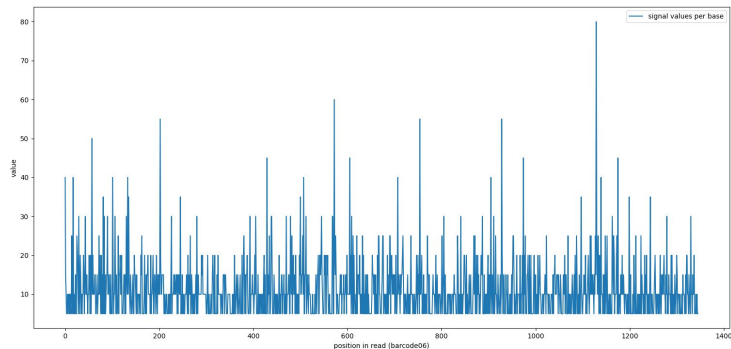


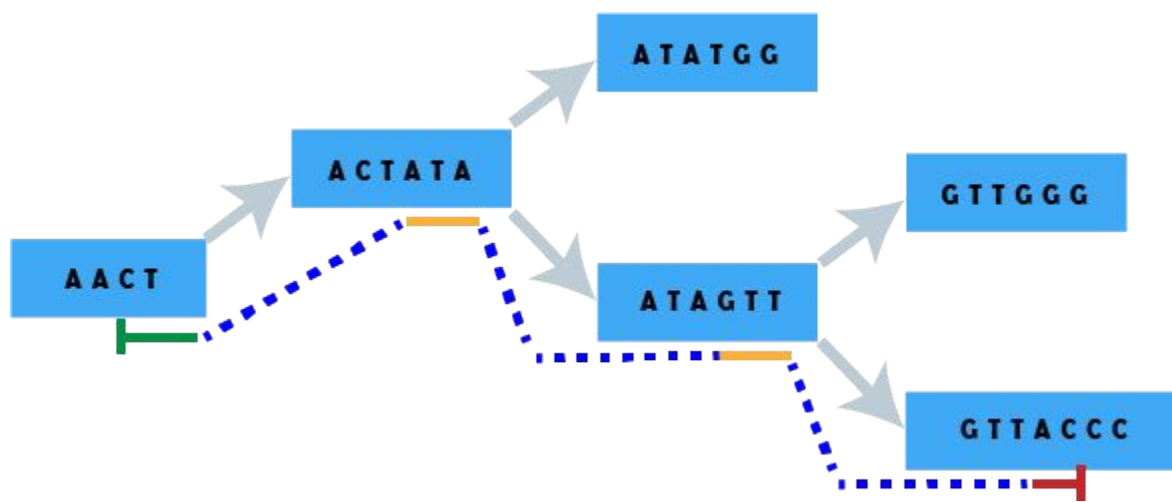
reprezentácia sekvencií v pamäti,
De Bruijnov graf,
hľadanie sledu,
aplikácia

II. Postup



Každý vrchol obsahuje “falošný” umelo vygenerovaný squiggle
 Pórový model (mean)
 Počet hodnôt per k-mer





“Hlava”, “telo”, “chvost” hľadanej cesty

Prechody grafom

Práca s maticou:

1. Pokrok v rámci vstupného readu
2. Lacnejšie hodnoty pre “hlavu”
3. Posledný riadok a stĺpec sú významné

pos	-	1	2	3	4	5
-	val	3	2	8	3	4
1	2	1	1	7	8	10
2	8	6	7	1	6	10
3	7	10	11	2	5	8
4	9	16	17	3	8	10
5	2	17	16	9	4	6

Hľadanie počiatočných vrcholov cesty

Skórovacia funkcia pre posúdenie, či by vrchol bol dobrou “hlavou”

Najlepšie 1% vrcholov (zaokrúhlené nahor)

...

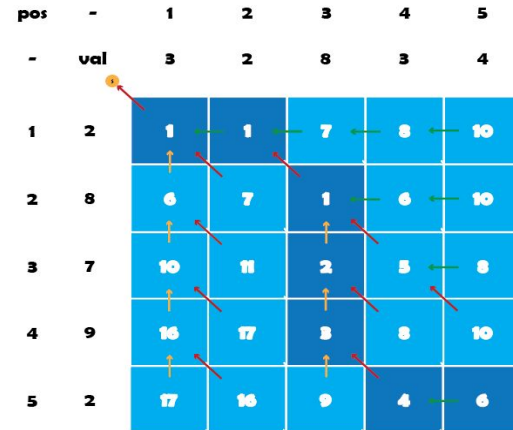
```
double read_aligned_multip = N * 3.0 / (read_position * 1.0);
```

```
double read_pos_fraction = read_position * 1.0 / N;
```

```
double score = (penalty+read_pos_fraction) * multip * read_aligned_multip *
```

```
read_aligned_multip;
```

```
return score;
```



Rozširovanie cesty

```
struct AlignResult {  
    double penalty = 0;  
    int read_pos = 0;  
    long long path_tag = -1;  
    int len = 0;  
};
```

DTW vo vnútri vrcholu - pole vstupných a výstupných
“situácií”

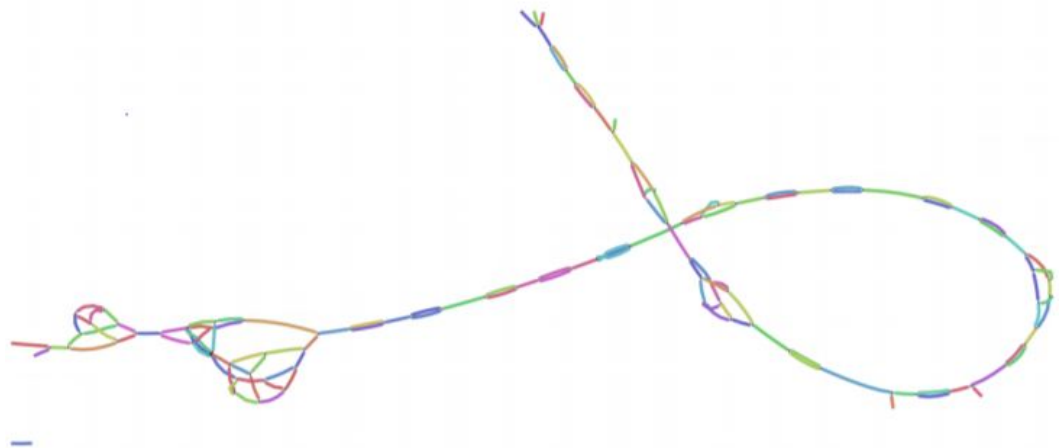
Kontrola celkovej “situácie”

```
(penalty*multip) / (read_position*1.0);
```

III. Experiment

	Dataset 1	Dataset 2
Vertices	148	159
Edges	200	255
Total sequence length	4209	2984

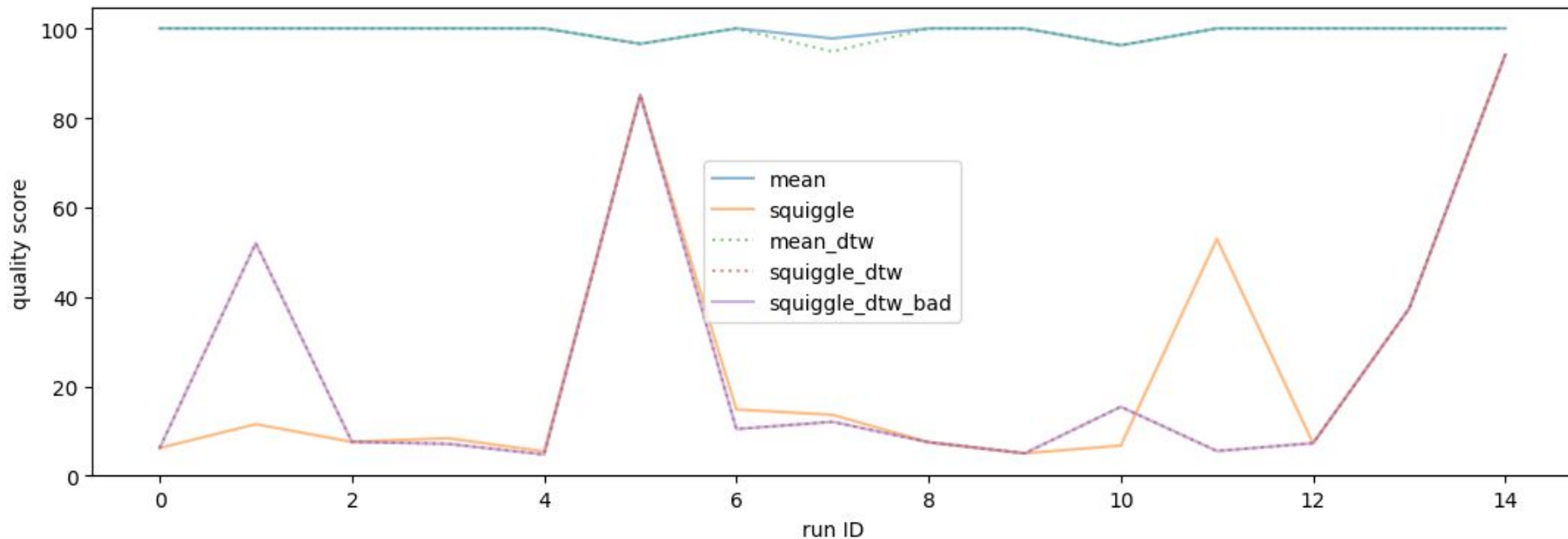
Fusobacterium nucleatum



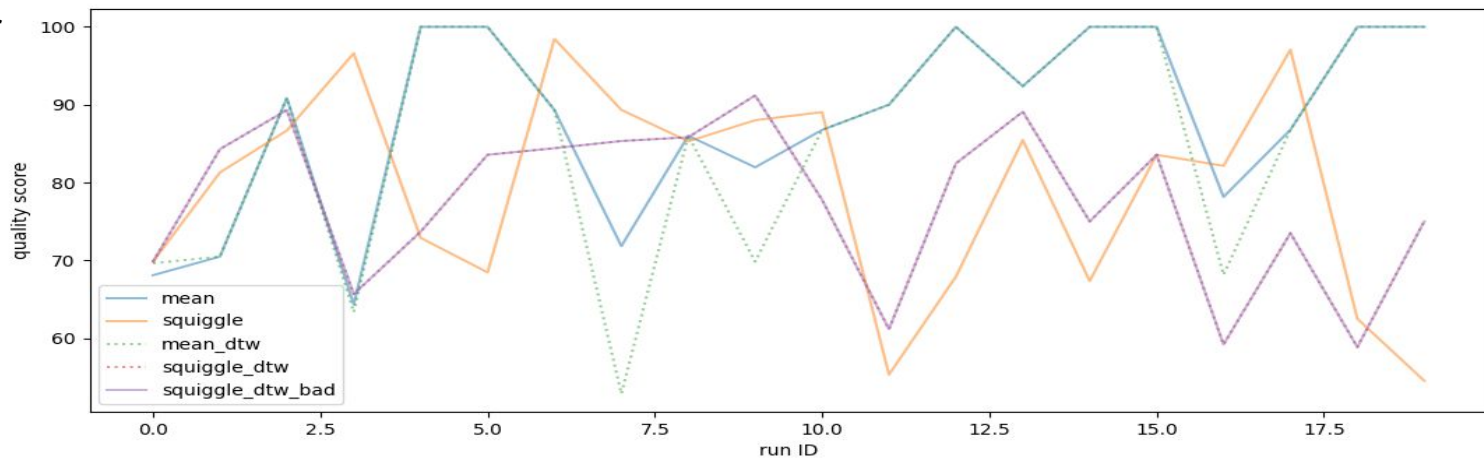
Generovanie input squiggle + referenčný squiggle

Cost function v DTW (eu, p2, pc)

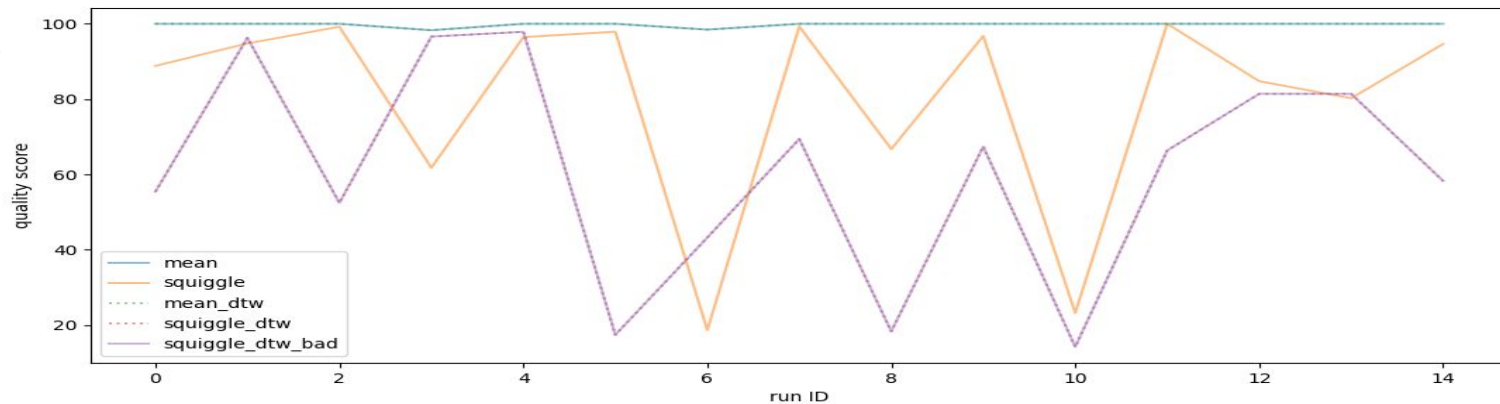
Úprava penalty podľa dĺžky warping path (yes, no)



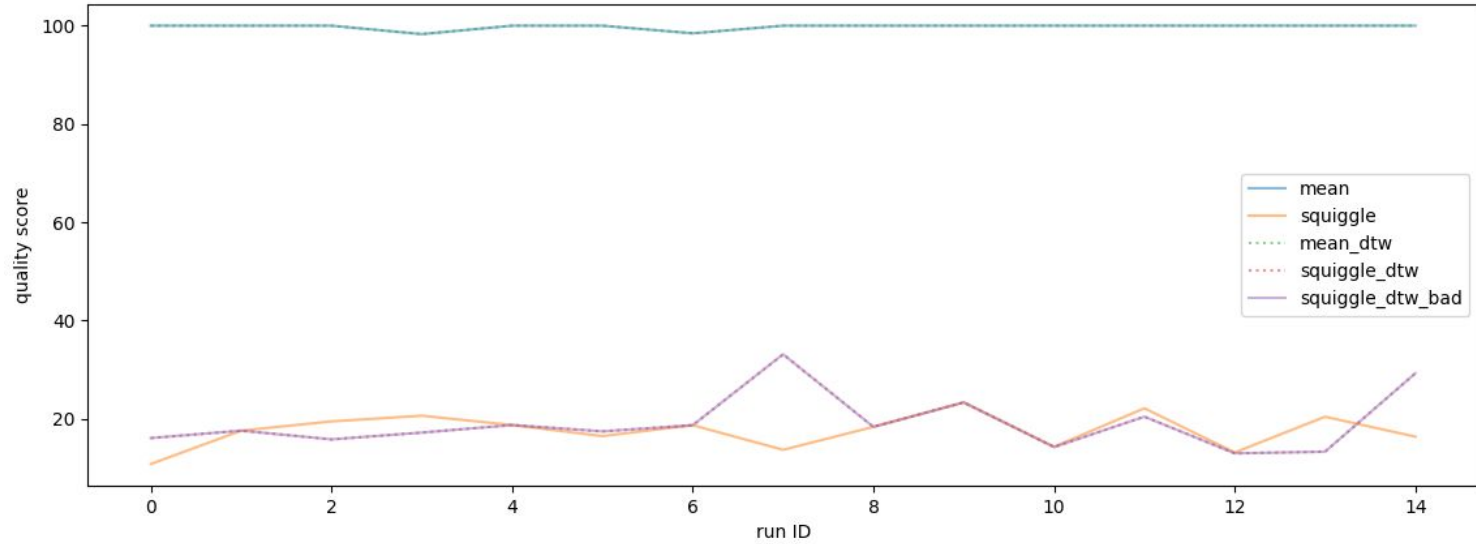
eu-l-r



p2-n-r



pc-n-r



Výkon aj kvalita sa líšia v závislosti od topológie grafu
pc a eu “zlé funkcie” - p2 trochu náprava - I viac náprava
Všetky squiggle rany majú podobné výsledky - počty
generovaných hodnôt v referencii nie sú podstatné (diskriminácia
všetkých vs. len niekoho)

IV. Súhrn a námety na ďalšie štúdium

Experimentálny postup, začať s vymysleným pekným readom,
zhoršovať ho
Výsledky sú zlé

Prečo je prvý dataset taký neúspešný?

Topológia skutočného grafu ??
Úprava vyhľadávania, zložitosť, veľké datasety

Veľa šumu

Ďakujem za pozornosť