

Detekcia tandemových opakovaní v nanopórových dátach

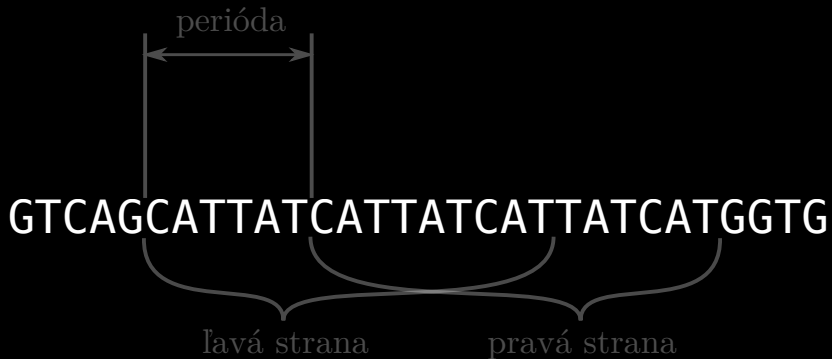
Eduard Batmendijn
Školiteľ: doc. Mgr. Tomáš Vinař, PhD.

18. júna 2020

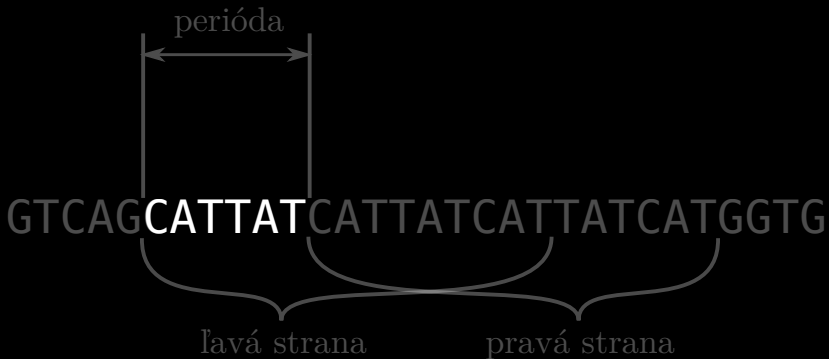
Detekcia tandemových opakovaní v nanopórových dátach

Detekcia tandemových
opakovaní v nanopórových
dátach

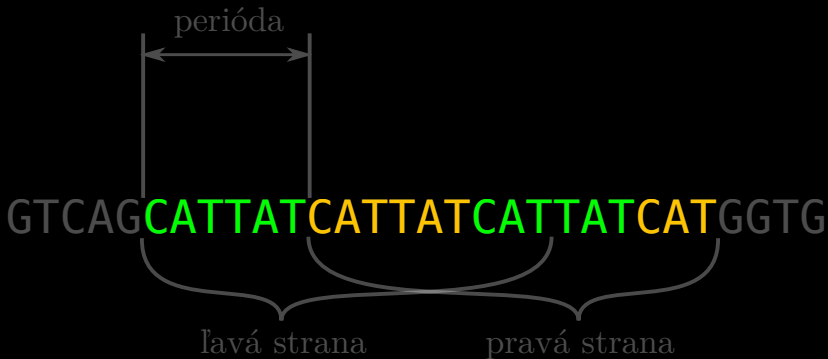
Tandemové opakovanie



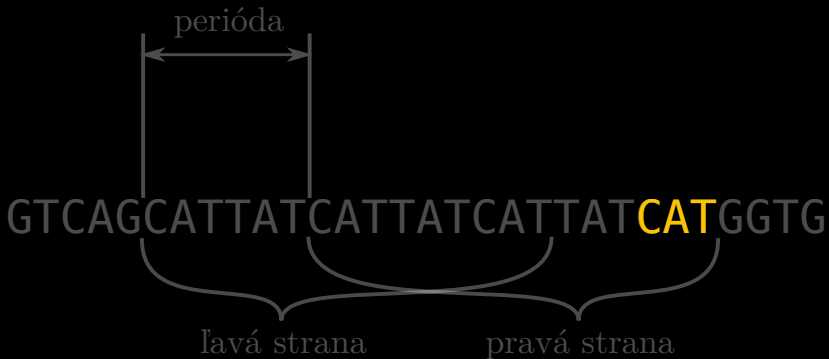
Tandemové opakovanie



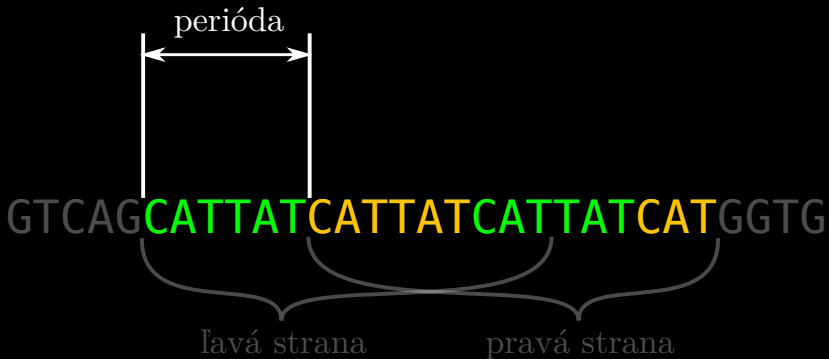
Tandemové opakovanie



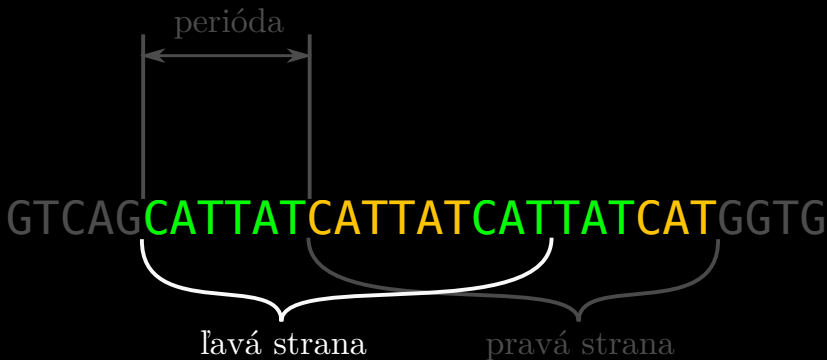
Tandemové opakovanie



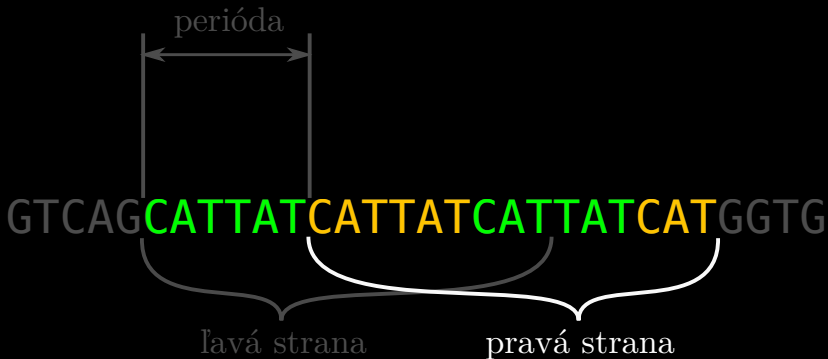
Tandemové opakovanie



Tandemové opakovanie



Tandemové opakovanie



Nedokonalé tandemové opakovanie

GTCAGCATTATCATGATCATTATCATGGTG

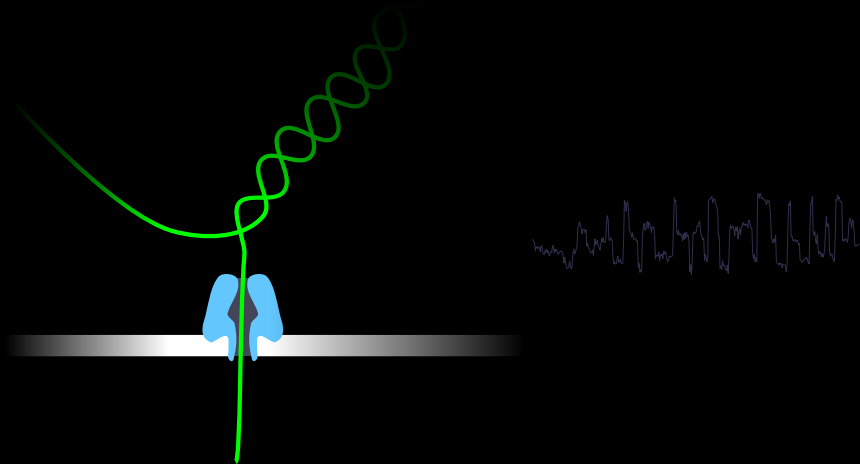
Nedokonalé tandemové opakovanie

GTCAGCATTATCATGATCATTATCATGGTG

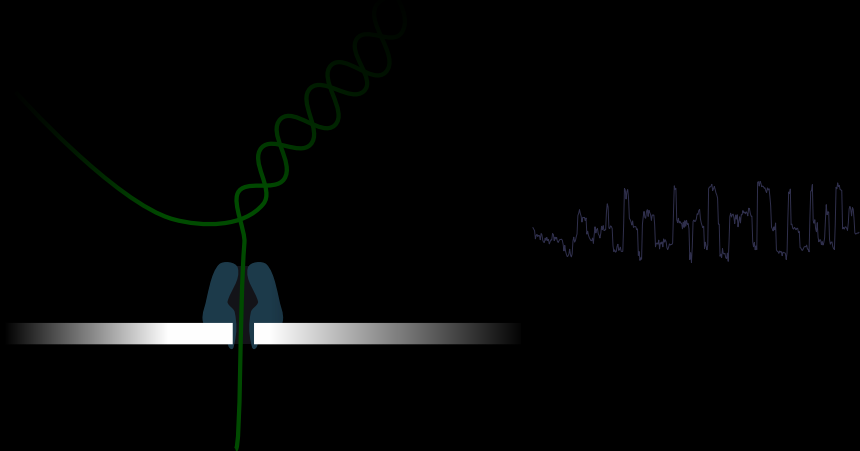
Detekcia tandemových opakovaní v nanopórových dátach

Detekcia tandemových
opakovaní v nanopórových
dátach

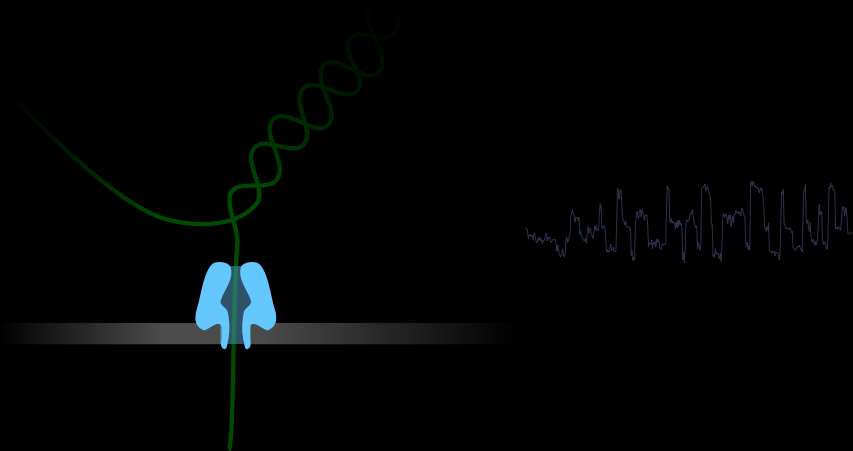
Nanopórové sekvenovanie



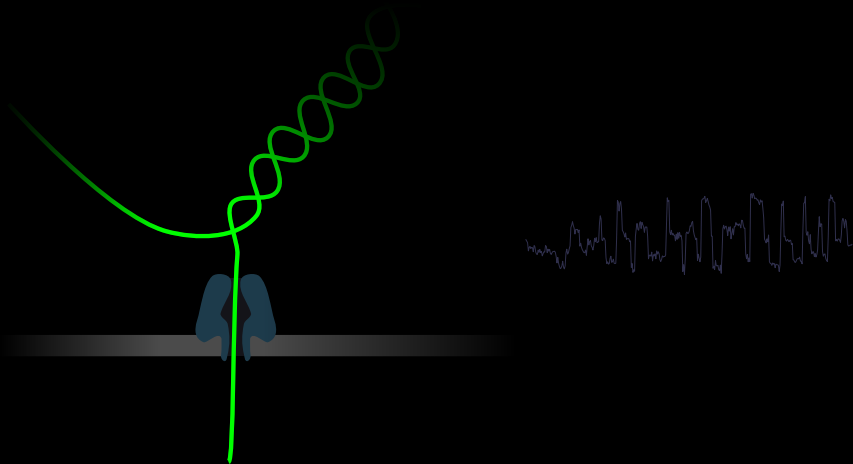
Nanopórové sekvenovanie



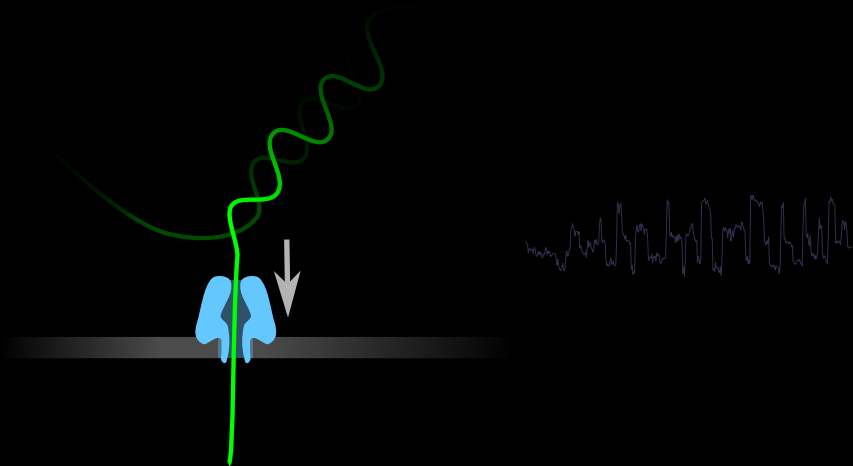
Nanopórové sekvenovanie



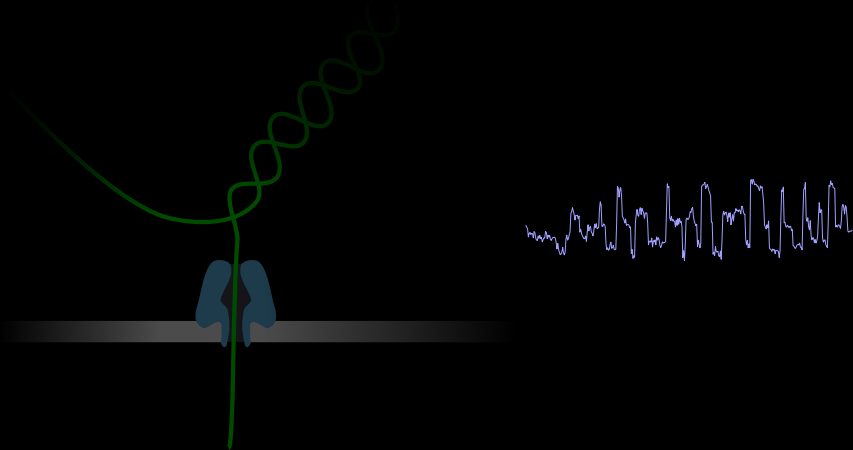
Nanopórové sekvenovanie



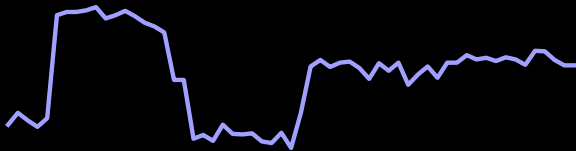
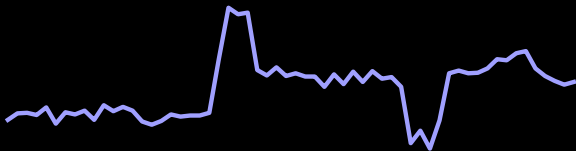
Nanopórové sekvenovanie



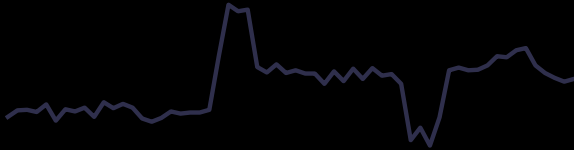
Nanopórové sekvenovanie



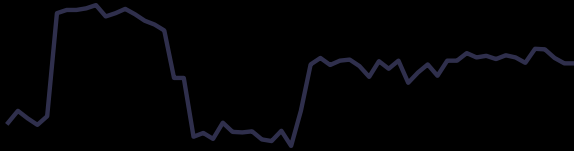
Nerovnomerná rýchlosť



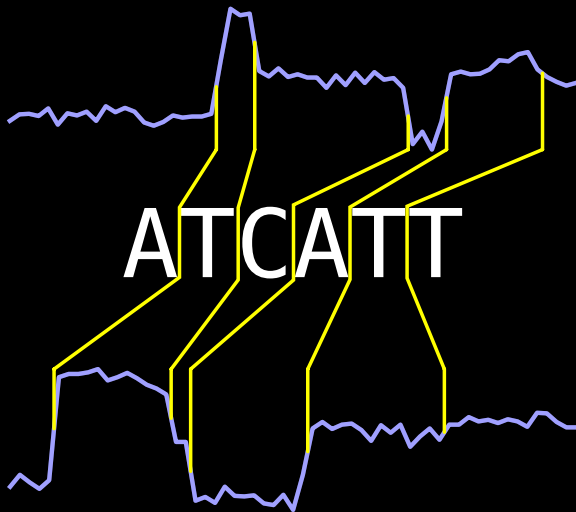
Nerovnomerná rýchlosť



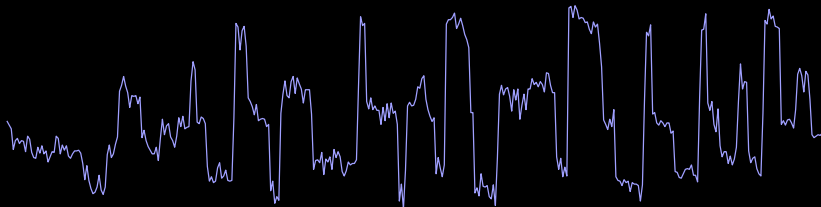
ATCATT



Nerovnomerná rýchlosť

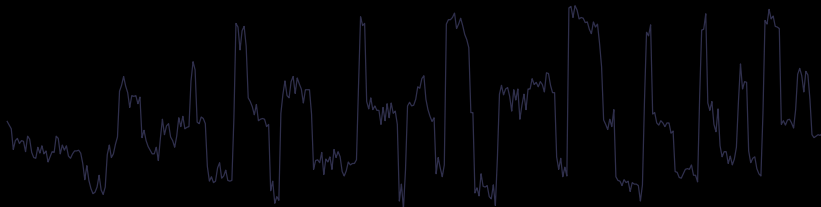


Určovanie báz (basecalling)



AAAAAAAAACACTTTCATCATTATCA
TTATCATTATCATCATCATATCTA

Určovanie báz (basecalling)

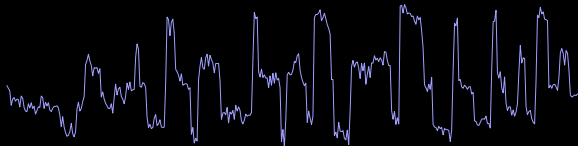


AAAAAAACA TTTTCATCATTATCA
TTATCATTATCATCATATCTA

85% – 90%

Riešenie

Naivná detekcia

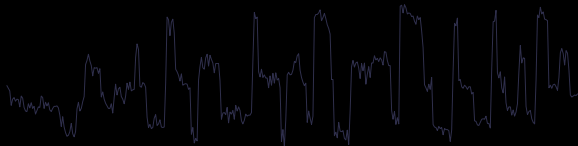


AAAAAAACACTTTCATCATTATCA
TTATCATTATCATCATATCTA



AAAAAAACACTTTCATCATTATCA
TTATCATTATCATCATATCTA

Naivná detekcia

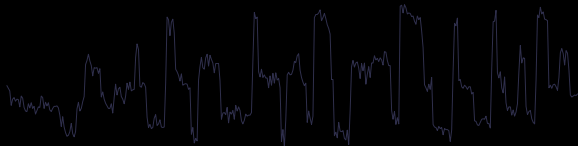


AAAAAAACACTTTCATCATTATCA
TTATCATTATCATCATATCTA



AAAAAAACACTTTCATCATTATCA
TTATCATTATCATCATATCTA

Naivná detekcia



AAAAAAACACTTTCATTATCA
TTATCATTATCATCATATCTA



AAAAAAACACTTTCATCATTATCA
TTATCATTATCATCATATCTA

Naivná detekcia

- ▶ Chybovosť prekladu
- ▶ Tvrdá predpoveď
- ▶ Nápad 1: detekcia priamo v signáli
- ▶ Nápad 2: čiastočné predspracovanie basecallerom

Naivná detekcia

- ▶ Chybovosť prekladu
- ▶ Tvrdá predpoveď
- ▶ Nápad 1: detekcia priamo v signáli
- ▶ Nápad 2: čiastočné predspracovanie basecallerom

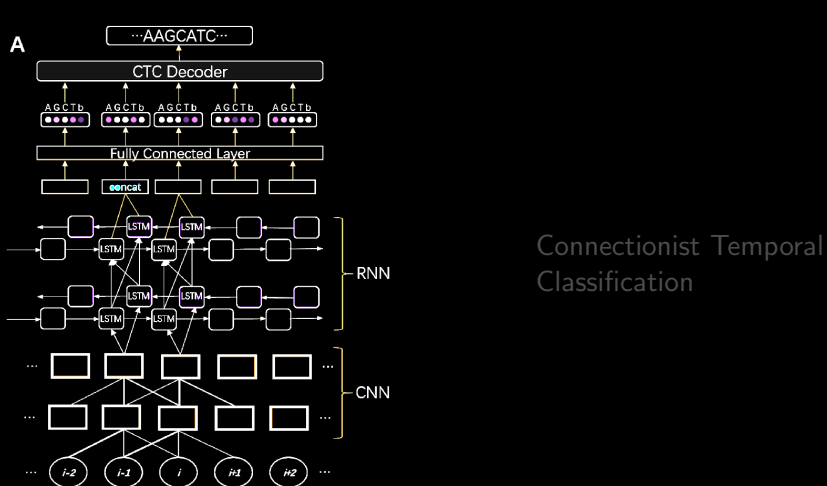
Naivná detekcia

- ▶ Chybovosť prekladu
- ▶ Tvrdá predpoveď
- ▶ **Nápad 1: detekcia priamo v signáli**
- ▶ Nápad 2: čiastočné predspracovanie basecallerom

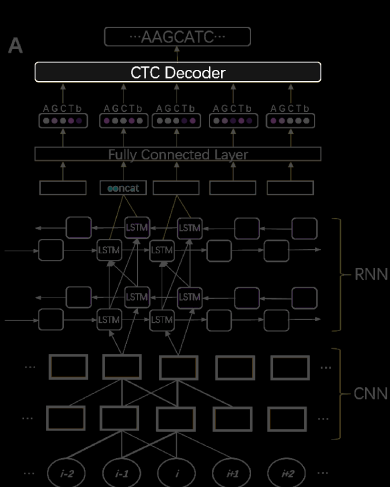
Naivná detekcia

- ▶ Chybovosť prekladu
- ▶ Tvrdá predpoveď
- ▶ Nápad 1: detekcia priamo v signáli
- ▶ Nápad 2: čiastočné predspracovanie basecallerom

Chiron (Teng et al. 2018)

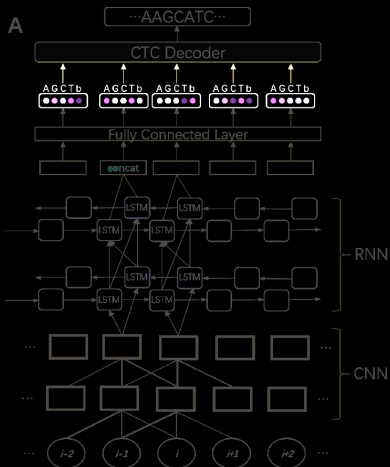


Chiron (Teng et al. 2018)



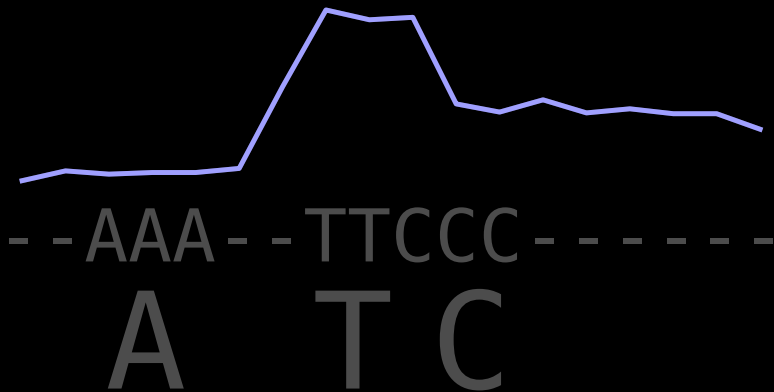
Connectionist Temporal
Classification

Chiron (Teng et al. 2018)



Connectionist Temporal
Classification

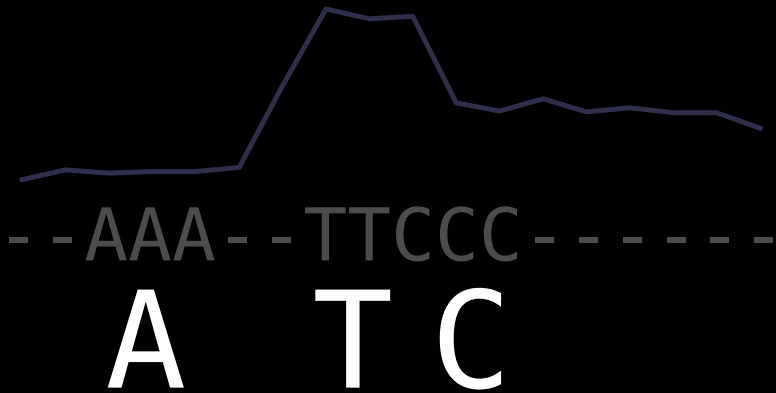
CTC anotácie



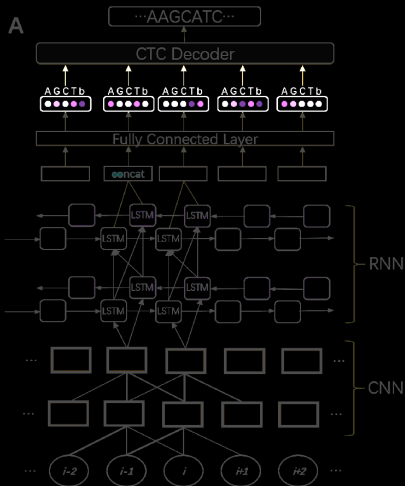
CTC anotácie



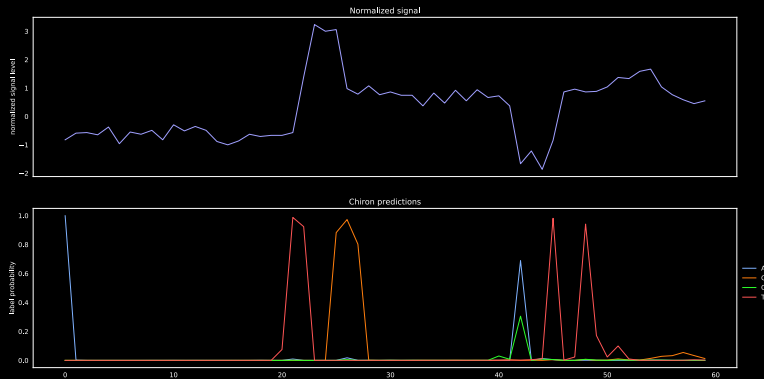
CTC anotácie



CTC predpovede



CTC predpovede



CTC dekodovanie

--AAA--TTCCC-----

A-----T-C-----

AAAAAATTTCCCCCCCCC

ATC-----

ATC

--AAA--GGGGG-----

-A-----G-----

AAAAAAGGGGGGGGGGGG

AG-----

AG

CTC dekodovanie

--AAA--TTCCC-----

A-----T-C-----

AAAAAATTTCCCCCCCCC

ATC-----

ATC

--AAA--GGGGG-----

-A-----G-----

AAAAAAGGGGGGGGGGGG

AG-----

AG

CTC dekodovanie

--AAA--TTCCC-----

A-----T-C-----

AAAAAATTTCCCCCCCCC

ATC-----

ATC

--AAA--GGGGG-----

-A-----G-----

AAAAAAGGGGGGGGGGGG

AG-----

AG

Detekcia v CTC predpovediach

- ▶ Tantan (Frith 2011)
- ▶ Skrytý Markovovský model (HMM)
- ▶ Rozšírenie z nukleotidov na CTC predpovede

Detekcia v CTC predpovediach

- ▶ Tantan (Frith 2011)
- ▶ Skrytý Markovovský model (HMM)
- ▶ Rozšírenie z nukleotidov na CTC predpovede

Detekcia v CTC predpovediach

- ▶ Tantan (Frith 2011)
- ▶ Skrytý Markovovský model (HMM)
- ▶ Rozšírenie z nukleotidov na CTC predpovede

Rozšiřovanie HMM

- ▶ HMM pre nukleotidy
- ▶ blanky
- ▶ nedokonalosti
- ▶ zlučovanie rovnakých označení

Rozširovanie HMM

- ▶ HMM pre nukleotidy
- ▶ **blanky**
- ▶ nedokonalosti
- ▶ zlučovanie rovnakých označení

Rozšiřovanie HMM

- ▶ HMM pre nukleotidy
- ▶ blanky
- ▶ **nedokonalosti**
- ▶ zlučovanie rovnakých označení

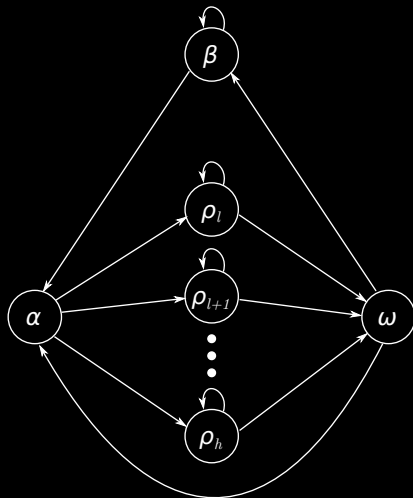
Rozšiřovanie HMM

- ▶ HMM pre nukleotidy
- ▶ blanky
- ▶ nedokonalosti
- ▶ zlučovanie rovnakých označení

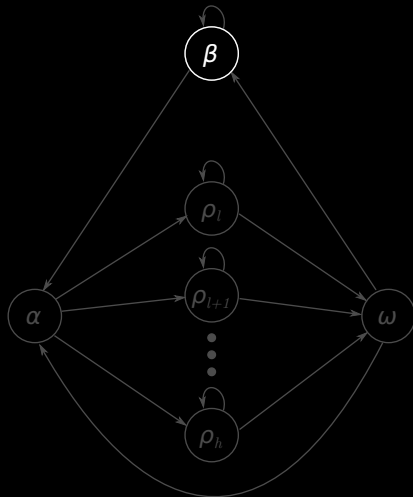
Rozšiřovanie HMM

- ▶ HMM pre nukleotidy
- ▶ blanky
- ▶ nedokonalosti
- ▶ zlučovanie rovnakých označení

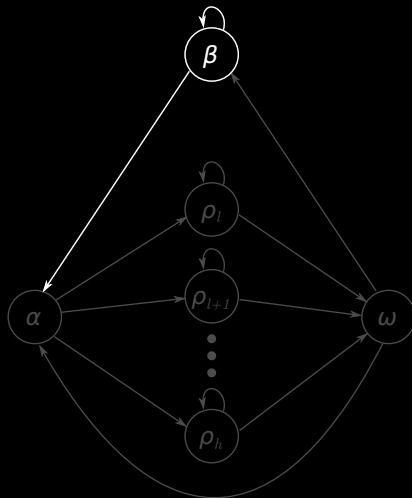
HMM pre nukleotidy



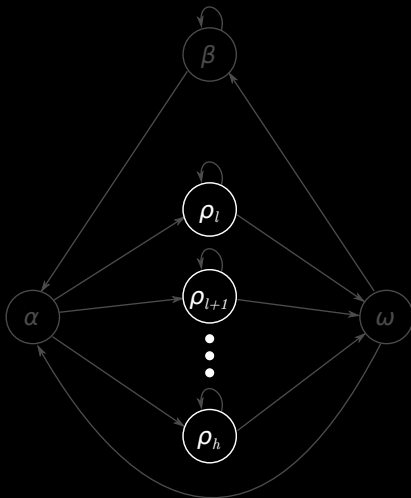
HMM pre nukleotidy



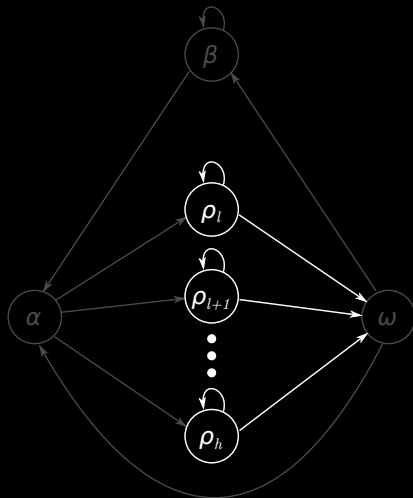
HMM pre nukleotidy



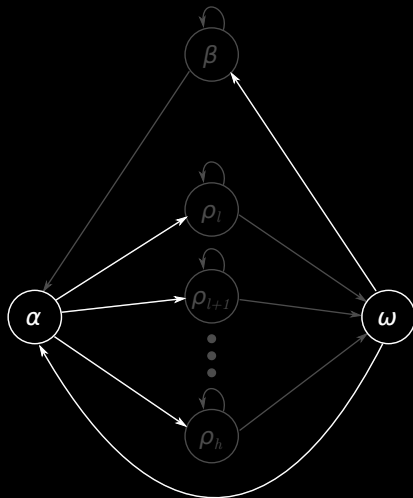
HMM pre nukleotidy



HMM pre nukleotidy

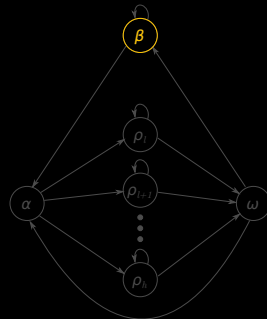


HMM pre nukleotidy



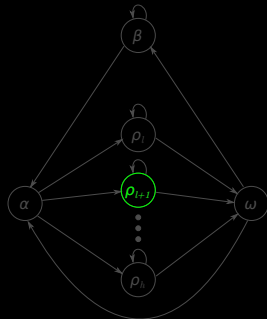
HMM pre nukleotidy

GGTCATCATCAA



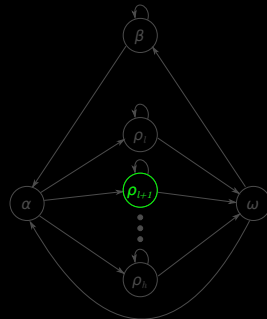
HMM pre nukleotidy

GGTCATCATCATCAA



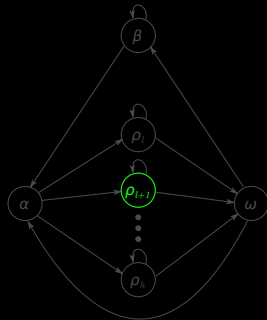
HMM pre nukleotidy

GGTCATCATCATCAA



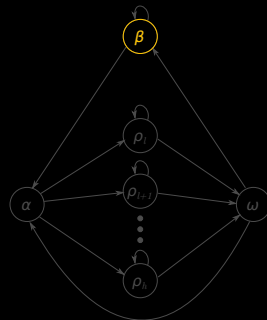
HMM pre nukleotidy

GGTCATCATCAA



HMM pre nukleotidy

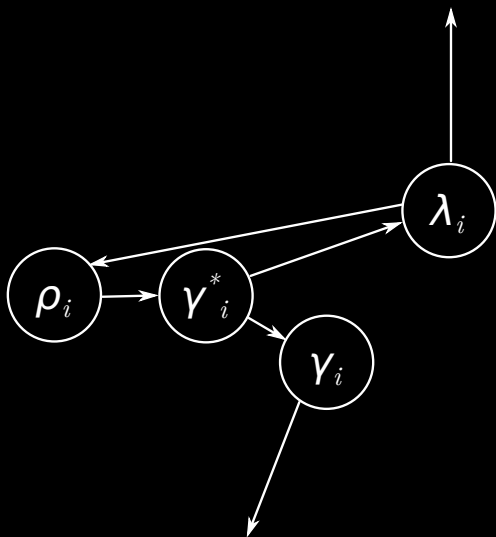
GGTCATCATCAA



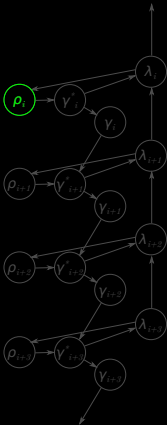
Rozširovanie HMM

- ▶ HMM pre nukleotidy
- ▶ **blanky**
- ▶ nedokonalosti
- ▶ zlučovanie rovnakých označení

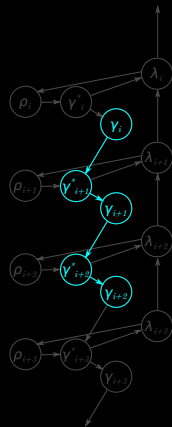
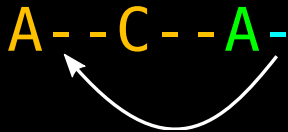
Blanky



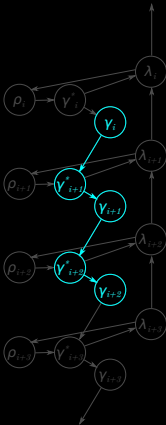
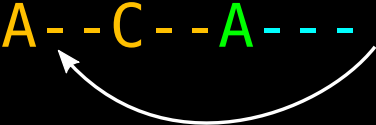
Blanky



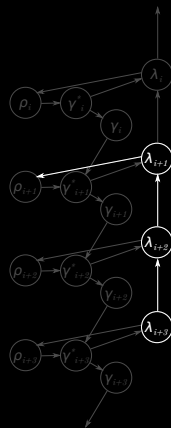
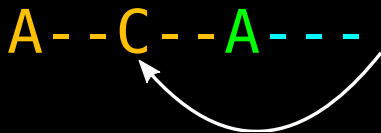
Blanky



Blanky



Blanky



Rozšiřovanie HMM

- ▶ HMM pre nukleotidy
- ▶ blanky
- ▶ **nedokonalosti**
- ▶ zlučovanie rovnakých označení

Nedokonalosti

- ▶ zmenené nukleotidy
- ▶ emisná distribúcia v ρ
- ▶ inzercie
- ▶ emisná distribúcia v γ
- ▶ delécie
- ▶ prechody z λ

Nedokonalosti

- ▶ zmenené nukleotidy
- ▶ emisná distribúcia v ρ
- ▶ inzercie
- ▶ emisná distribúcia v γ
- ▶ delécie
- ▶ prechody z λ

Nedokonalosti

- ▶ zmenené nukleotidy
- ▶ emisná distribúcia v ρ
- ▶ **inzercie**
- ▶ emisná distribúcia v γ
- ▶ delécie
- ▶ prechody z λ

Nedokonalosti

- ▶ zmenené nukleotidy
- ▶ emisná distribúcia v ρ
- ▶ inzercie
- ▶ emisná distribúcia v γ
- ▶ delécie
- ▶ prechody z λ

Nedokonalosti

- ▶ zmenené nukleotidy
- ▶ emisná distribúcia v ρ
- ▶ inzercie
- ▶ emisná distribúcia v γ
- ▶ **delécie**
- ▶ prechody z λ

Nedokonalosti

- ▶ zmenené nukleotidy
- ▶ emisná distribúcia v ρ
- ▶ inzercie
- ▶ emisná distribúcia v γ
- ▶ delécie
- ▶ prechody z λ

Rozširovanie HMM

- ▶ HMM pre nukleotidy
- ▶ blanky
- ▶ nedokonalosti
- ▶ zlučovanie rovnakých označení

Zlučovanie rovnakých označení

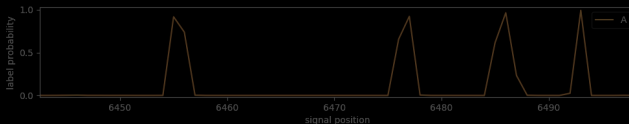
- ▶ modelovanie: zhruba 15-krát viac stavov
- ▶ ignorovanie: nižšia presnosť



- ▶
- ▶ Heuristika: lokálne maximá

Zlučovanie rovnakých označení

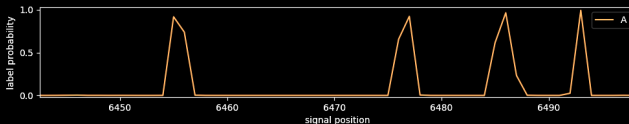
- ▶ modelovanie: zhruba 15-krát viac stavov
- ▶ ignorovanie: nižšia presnosť



- ▶
- ▶ Heuristika: lokálne maximá

Zlučovanie rovnakých označení

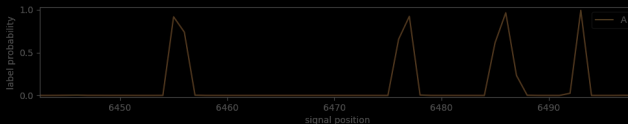
- ▶ modelovanie: zhruba 15-krát viac stavov
- ▶ ignorovanie: nižšia presnosť



- ▶
- ▶ Heuristika: lokálne maximá

Zlučovanie rovnakých označení

- ▶ modelovanie: zhruba 15-krát viac stavov
- ▶ ignorovanie: nižšia presnosť



- ▶ Heuristika: lokálne maximá

Vyhodnotenie

Datasety

Name	Reads	Repeats ratio	Average repeat period
Saping	100	7.61%	5.84
Jamang	100	18.07%	63.72

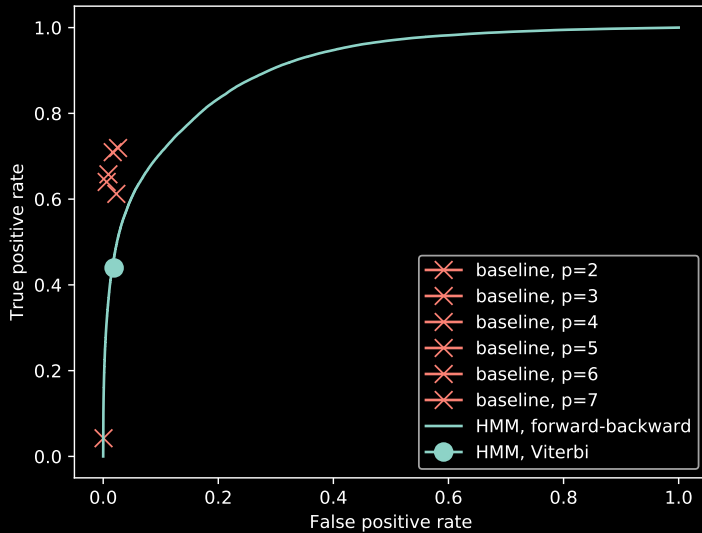
Metriky

- ▶ True positive rate
- ▶ False positive rate

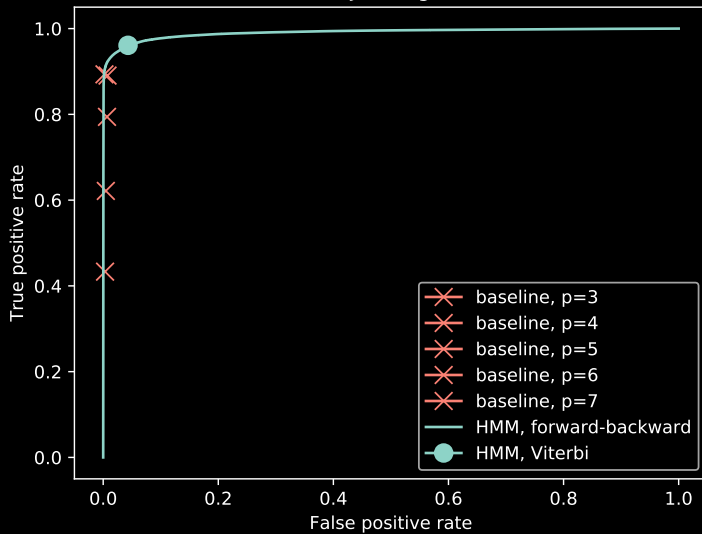
Metriky

- ▶ True positive rate
- ▶ False positive rate

Saping



Jamang



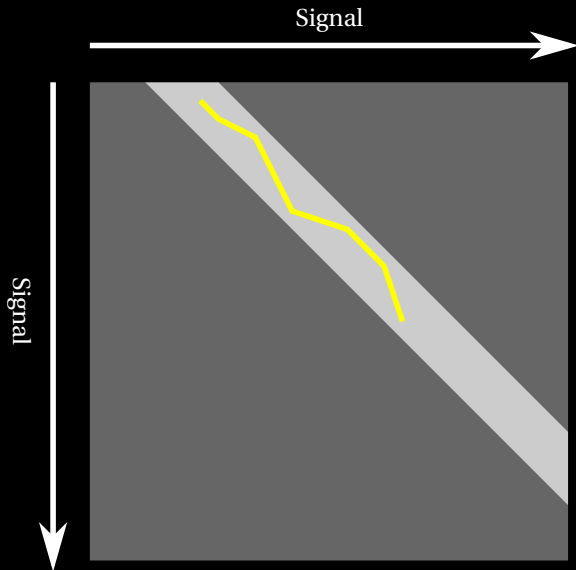
Použité obrázky

H. Teng et al. Chiron: translating nanopore raw signal directly into nucleotide sequence using deep learning. *GigaScience*, 7(5), 04 2018. [giy037](https://doi.org/10.1093/gigascience/giy037).

Otázka 1

- ▶ Akú funkciu optimalizuje algoritmus z kapitoly 3?
- ▶ Dala by sa upraviť tak, aby bolo možné použiť efektívnejší algoritmus?

Otázka 1



Otázka 1

- ▶ Akú funkciu optimalizuje algoritmus z kapitoly 3?
- ▶ Dala by sa upraviť tak, aby bolo možné použiť efektívnejší algoritmus?

Otázka 1

- ▶ $s_1 s_2 \dots s_n$
- ▶ $t_1 t_2 \dots t_m$
- ▶ $s_{i_1} s_{i_2} \dots s_{i_l}$
- ▶ $t_{j_1} t_{j_2} \dots t_{j_l}$
- ▶ $i_k \leq i_{k+1} \leq i_k + 1$
- ▶ $l \geq m, n$
- ▶ $\frac{b(n+m)}{2} - \sum_{k=1}^l (s_{i_k} - t_{i_k})^2$
- ▶ $b = \text{bonus}$

Otázka 1

▶ $s_1 s_2 \dots s_n$

▶ $t_1 t_2 \dots t_m$

▶ $s_{i_1} s_{i_2} \dots s_{i_l}$

▶ $t_{j_1} t_{j_2} \dots t_{j_l}$

▶ $i_k \leq i_{k+1} \leq i_k + 1$

▶ $l \geq m, n$

▶ $\frac{b(n+m)}{2} - \sum_{k=1}^l (s_{i_k} - t_{i_k})^2$

▶ $b = \text{bonus}$

Otázka 1

- ▶ $s_1 s_2 \dots s_n$
- ▶ $t_1 t_2 \dots t_m$
- ▶ $s_{i_1} s_{i_2} \dots s_{i_l}$
- ▶ $t_{j_1} t_{j_2} \dots t_{j_l}$
- ▶ $i_k \leq i_{k+1} \leq i_k + 1$
- ▶ $l \geq m, n$
- ▶ $\frac{b(n+m)}{2} - \sum_{k=1}^l (s_{i_k} - t_{i_k})^2$
- ▶ $b = \text{bonus}$

Otázka 1

▶ $s_1 s_2 \dots s_n$

▶ $t_1 t_2 \dots t_m$

▶ $s_{i_1} s_{i_2} \dots s_{i_l}$

▶ $t_{j_1} t_{j_2} \dots t_{j_l}$

▶ $i_k \leq i_{k+1} \leq i_k + 1$

▶ $l \geq m, n$

▶ $\frac{b(n+m)}{2} - \sum_{k=1}^l (s_{i_k} - t_{i_k})^2$

▶ $b = \text{bonus}$

Otázka 1

▶ $s_1 s_2 \dots s_n$

▶ $t_1 t_2 \dots t_m$

▶ $s_{i_1} s_{i_2} \dots s_{i_l}$

▶ $t_{j_1} t_{j_2} \dots t_{j_l}$

▶ $i_k \leq i_{k+1} \leq i_k + 1$

▶ $l \geq m, n$

▶ $\frac{b(n+m)}{2} - \sum_{k=1}^l (s_{i_k} - t_{i_k})^2$

▶ $b = \text{bonus}$

Otázka 1

▶ $s_1 s_2 \dots s_n$

▶ $t_1 t_2 \dots t_m$

▶ $s_{i_1} s_{i_2} \dots s_{i_l}$

▶ $t_{j_1} t_{j_2} \dots t_{j_l}$

▶ $i_k \leq i_{k+1} \leq i_k + 1$

▶ $l \geq m, n$

▶ $\frac{b(n+m)}{2} - \sum_{k=1}^l (s_{i_k} - t_{i_k})^2$

▶ $b = \text{bonus}$

Otázka 1

▶ $s_1 s_2 \dots s_n$

▶ $t_1 t_2 \dots t_m$

▶ $s_{i_1} s_{i_2} \dots s_{i_l}$

▶ $t_{j_1} t_{j_2} \dots t_{j_l}$

▶ $i_k \leq i_{k+1} \leq i_k + 1$

▶ $l \geq m, n$

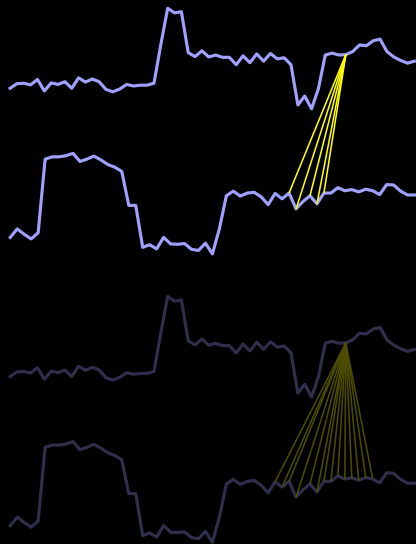
▶ $\frac{b(n+m)}{2} - \sum_{k=1}^l (s_{i_k} - t_{i_k})^2$

▶ $b = \text{bonus}$

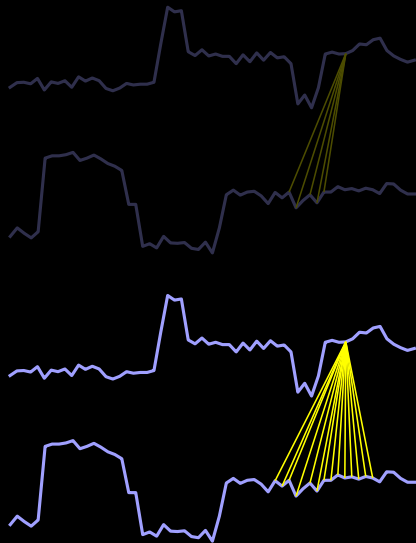
Otázka 1

- ▶ Akú funkciu optimalizuje algoritmus z kapitoly 3?
- ▶ Dala by sa upraviť tak, aby bolo možné použiť efektívnejší algoritmus?

Otázka 1



Otázka 1



Otázka 2

- ▶ Alternatívny prístup z konca časti 4.2.5 s použitím preprocessovania a jednoduchého modelu nebol až tak horší ako plný model. Vedeli by ste si predstaviť presun ďalších častí plného modelu do preprocessingu? Ak áno, tak ktoré?
- ▶ Dal by sa tiež použiť algoritmus z kapitoly 3?

Otázka 2

- ▶ Alternatívny prístup z konca časti 4.2.5 s použitím preprocessovania a jednoduchého modelu nebol až tak horší ako plný model. Vedeli by ste si predstaviť presun ďalších častí plného modelu do preprocessingu? Ak áno, tak ktoré?
- ▶ Dal by sa tiež použiť algoritmus z kapitoly 3?