

Využitie reťazcových grafov pri hľadani chýb v genómoch

Angelika Fedáková

Školiteľ: doc. Mgr. Bronislava Brejová, PhD.

01.07.2020

DNA

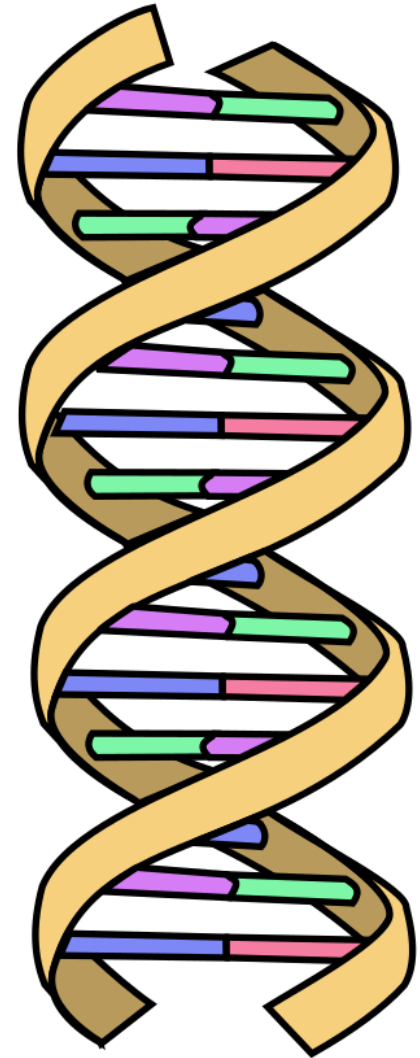
- ▶ Reť azec báz A,C,G,T
- ▶ Dve komplementárne vlákna

Chromozóm

- ▶ Súvislý úsek DNA

Genóm

- ▶ Súbor chromozómov



DNA

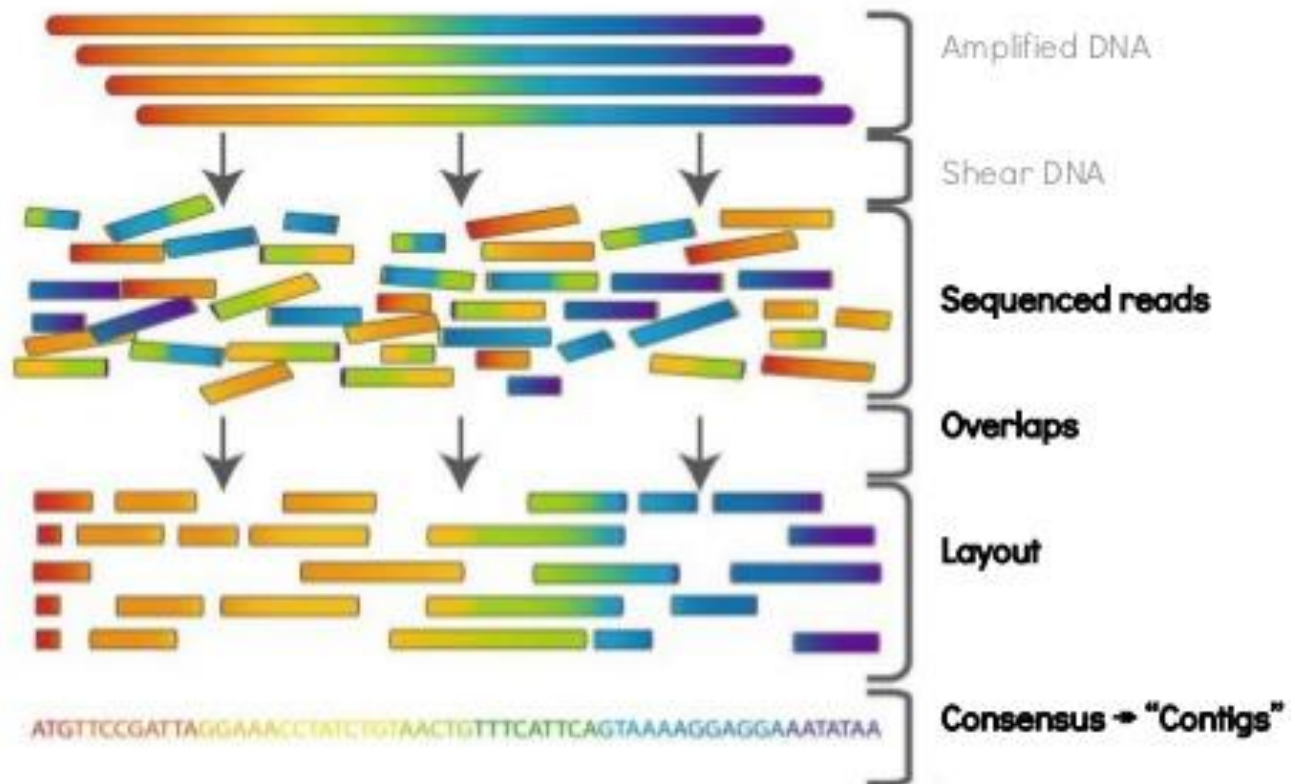
Sekvenovanie DNA

- ▶ Zisťovania báz chromozómov
- ▶ Sekvenovanie krátkych kúskov – čítaní
 - Krátke čítania, dlhé čítania

Zostavovanie genómu

- ▶ Proces skladania čítaní do pôvodných chromozómov

Overlap - Layout - Consensus



Zarovňavanie sekvencií

- ▶ Hľadanie podobností a rozdielov sekvencií
- ▶ Mapovanie báz sekvencií
- ▶ Mapovanie sekvencie ku grafu

Sekvencia 1: AGCTGGCTT

Sekvencia 2: GCAGGTCTT

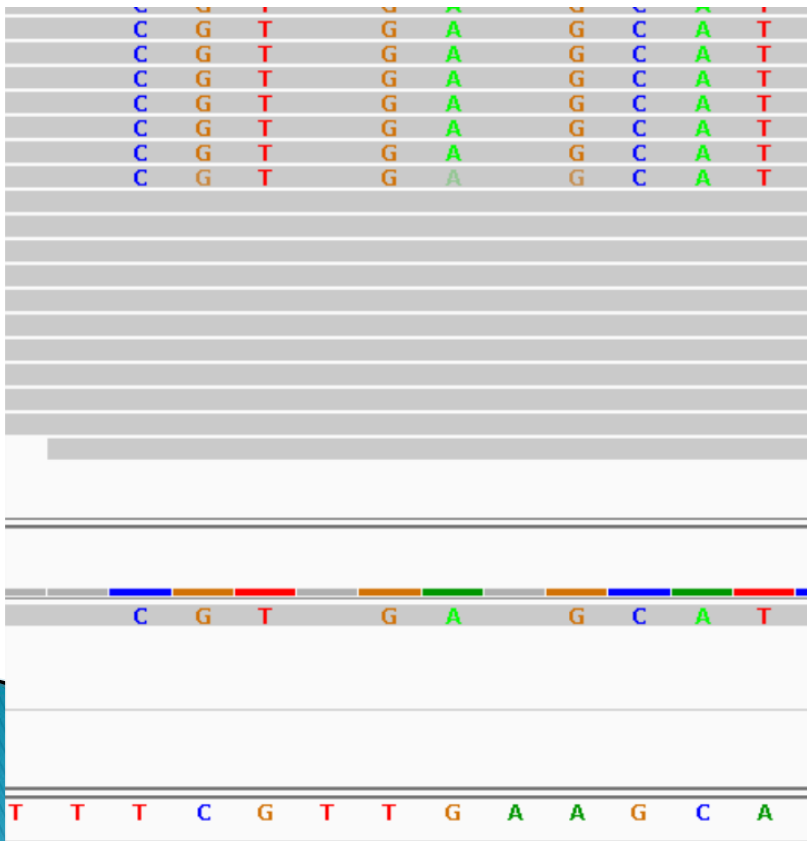
A	G	C	T	G	G	-	C	T	T
-	G	C	A	G	G	T	C	T	T
-1	+1	+1	-1	+1	+1	-1	+1	+1	+1

Skóre zarovnaní 4

Chyby DNA genómov

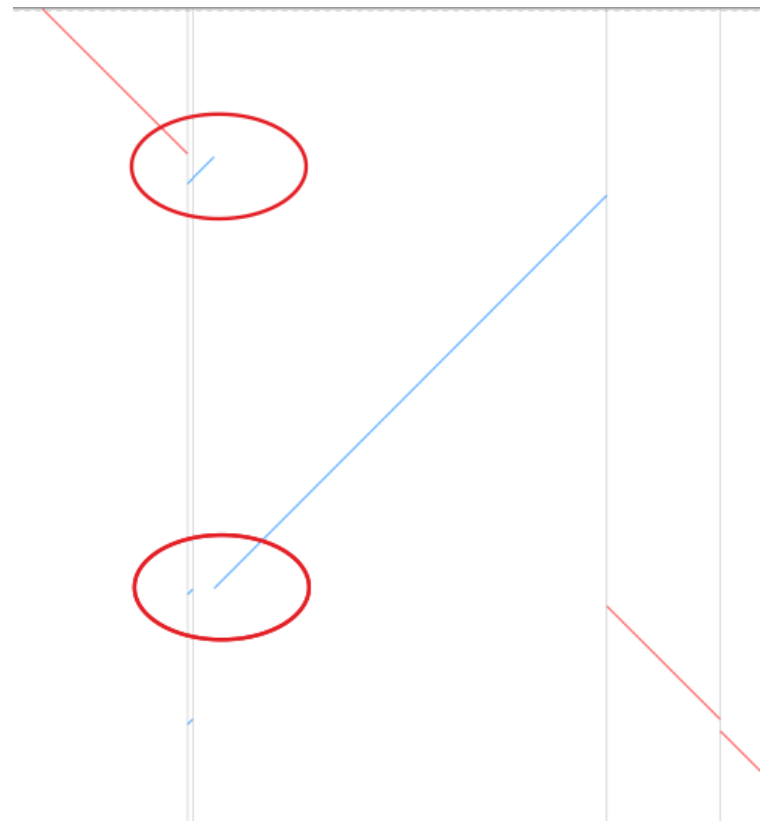
▶ Menšie chyby

- Zlé určenie jednotlivých báz

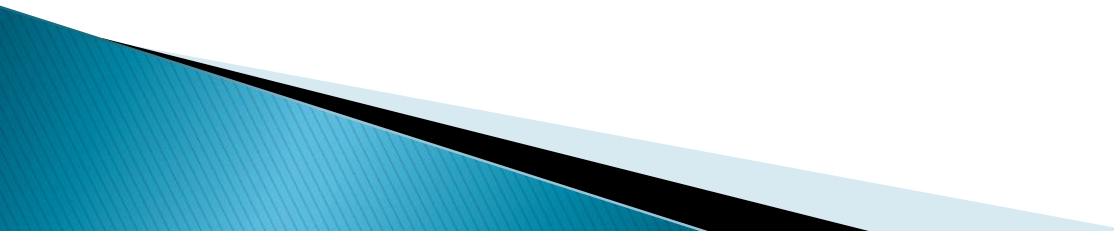


▶ Rozsiahlejšie chyby

- Zlé určenie poradia väčších kusov chromozómu



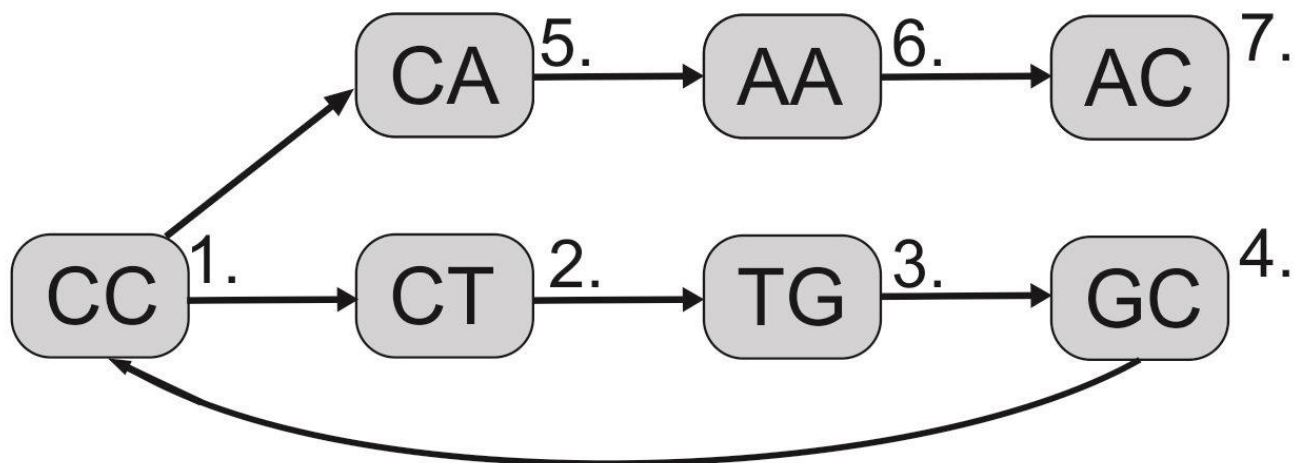
Ciele práce

- ▶ Priniest' nový pohľad na problém hľadania chýb
 - ▶ Identifikovať miesta s pravdepodobnosťou chyby
 - ▶ Porovnať výsledky s existujúcim softvérom
- 

Postup hľadania chýb

► Vytvorenie grafu

$k = 2$, krátke čítania: CCTG, TGCC, GCCA, CAAC



Sekvencia: TGCCAAC

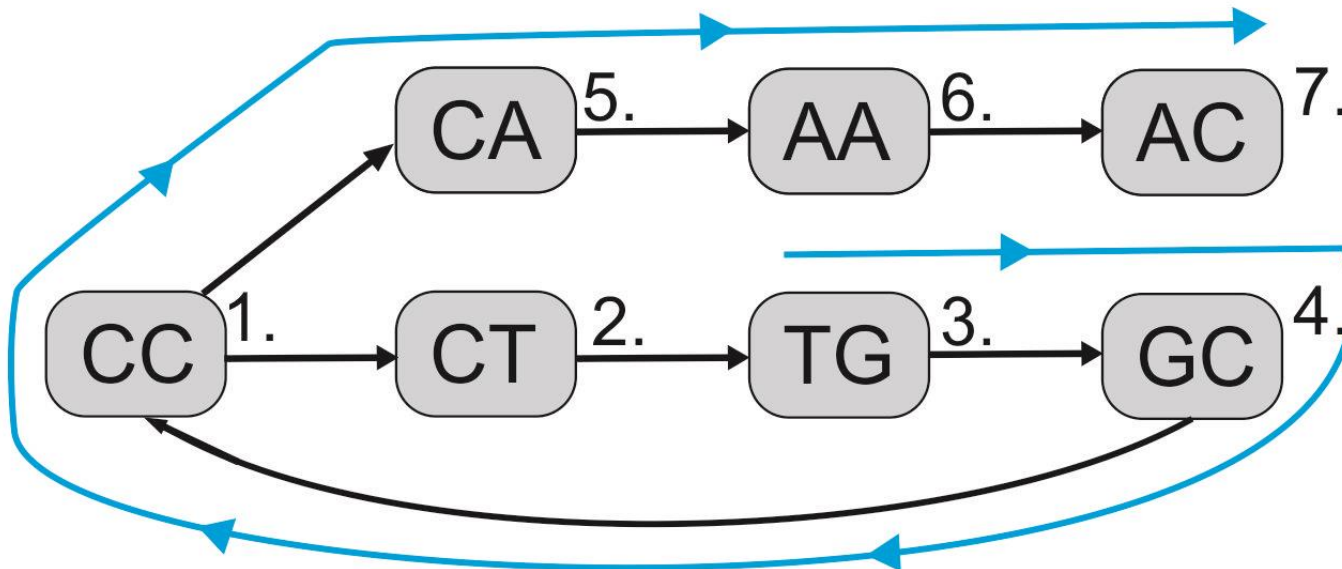
Dlhé čítania: CCTGCC, GCCAAC

Postup hľadania chýb

- ▶ Zarovnanie dlhých čítaní a sekvencie

Sekvencia: TGCCAAC

Dlhé čítania: CCTGCC, GCCAAC



Zarovnanie čítaní: 1. 2. 3. 4. 1. 4. 1. 5. 6. 7.

Zarovnanie sekvencie: 3. 4. 1. 5. 6. 7.

Vstup a výstup programu

- ▶ Vstup:
 - Skúmaný genóm/sekvencia (FASTA)
 - Krátke čítania (FASTA, FASTQ)
 - Dlhé čítania (FASTA, FASTQ)
- ▶ Výstup
 - Oblasti chybových miest (BED)

Experimentálne dáta

- ▶ Baktéria *Escherichia Coli* (E.~Coli)
 - Referenčná sekvencia (~500 000bp)
 - 5 zostavených sekvencií
- ▶ Kvasinka *Saprochaete fungicola*
 - Referenčná sekvencia (~20 000 000bp)
 - 2 zostavené sekvencie

Vyhodnotenie algoritmu

- ▶ Referenčné chyby, nájdené chyby, správne chyby
- ▶ Presnosť : $\frac{\text{správne nájdené chyby}}{\text{všetky nájdené chyby}}$
- ▶ Úplnosť : $\frac{\text{nájdené referenčné chyby}}{\text{všetky referenčné chyby}}$

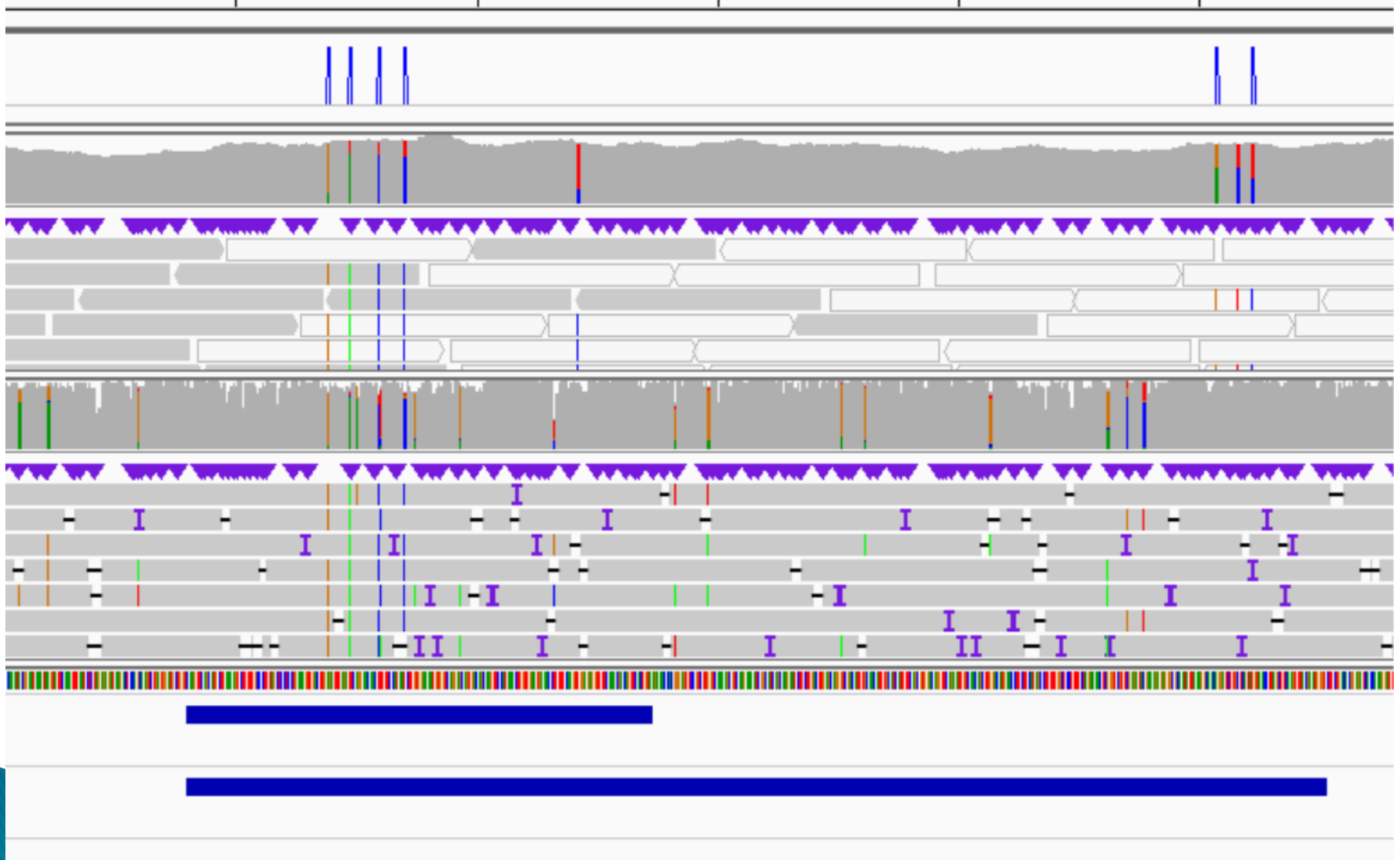
Prvé výsledky kvasinky

	<i>Spades</i>		<i>Miniasm</i>	
<i>Presnosť</i>				
<i>Limit1</i>	543/837	64.8%	325/480	67.7%
<i>Limit3</i>	1209/2591	46.6%	472/816	57.8%
<i>Limit5</i>	1310/3144	41.6%	593/1062	55.8%
<i>Limit7</i>	1379/3902	35.3%	662/1244	53.3%
<i>Úplnosť</i>				
<i>Limit1</i>	599/1599	37.4%	386/2513	15.3%
<i>Limit3</i>	1298/1599	81.2%	567/2513	22.5%
<i>Limit5</i>	1389/1599	86.8%	692/2513	27.5%
<i>Limit7</i>	1459/1599	91.2%	770/2513	30.6%

Finálne výsledky po vykonaných filtráciách

	<i>Spades</i>		<i>Miniasm</i>	
<i>Presnosť</i>				
<i>Limit1</i>	153/198	77.3%	101/136	74.3%
<i>Limit3</i>	439/615	71.4%	124/183	67.8%
<i>Limit5</i>	498/883	56.4%	153/227	67.4%
<i>Limit7</i>	544/1418	38.4%	174/263	66.2%
<i>Úplnosť</i>				
<i>Limit1</i>	185/661	28%	124/1506	8.2%
<i>Limit3</i>	459/661	69.4%	157/1506	10.4%
<i>Limit5</i>	515/661	77.9%	190/1506	12.6%
<i>Limit7</i>	546/661	82.5%	211/1506	14%

Príklad správne nájdenej chyby



Porovnanie so softvérom Pilon

	<i>Spades</i>	<i>Miniasm</i>
<i>Celá sekvencia</i>	0.39%	1%
<i>Referenčné chyby</i>	3.7%	38%
<i>Limit1</i>	5%	19.6%
<i>Limit3</i>	2.4%	8.4%
<i>Limit5</i>	2.5%	7.5%
<i>Limit7</i>	2.4%	8.1%

Možnosti d'alšej práce

- ▶ Zlepšenie zarovnaní ku grafu
- ▶ Objavenie d'alších filtrácií chýb
- ▶ Zložitejšia filtrácia homopolymérov
- ▶ Testovanie d'alších dát

Ďakujem za pozornosť



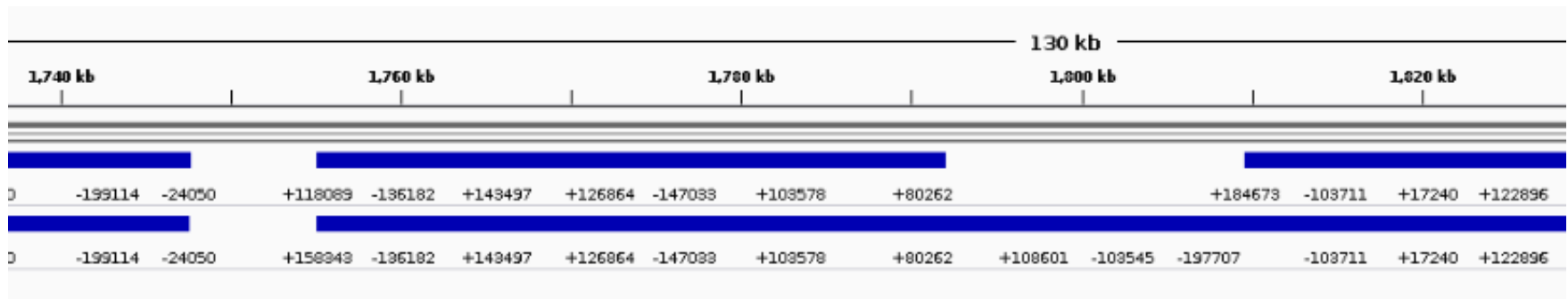
1. Ako veľmi by sa zmenili výsledky, ak by sme uvažovali realistickejšie prekryvy (napr. 50%)?

	<i>Spades</i>		<i>Miniasm</i>	
<i>Presnosť</i>				
<i>Prekryv</i>	<i>1bp</i>	<i>50%</i>	<i>1bp</i>	<i>50%</i>
<i>Limit1</i>	77.3%	65.2%	74.3%	55.1%
<i>Limit3</i>	71.4%	65.5%	67.8%	47.5%
<i>Limit5</i>	56.4%	51.9%	67.4%	51.5%
<i>Limit7</i>	38.4%	35%	66.2%	51%
<i>Úplnosť</i>				
<i>Prekryv</i>	<i>1bp</i>	<i>50%</i>	<i>1bp</i>	<i>50%</i>
<i>Limit1</i>	28%	24.1%	8.2%	6.8%
<i>Limit3</i>	69.4%	65.8%	10.4%	9%
<i>Limit5</i>	77.9%	75.2%	12.6%	11%
<i>Limit7</i>	82.5%	77.3%	14%	12.7%

2. Čím si vysvetľujete nekonzistentnosť v zarovnaní pôvodného a reverzného vlákna programom GraphAligner? Je k tomuto niečo v dostupnej literatúre, popr. kontaktovali ste tvorcov programu?

- ▶ Problém s nezarovnávaním rozsiahlejších kusov sekvencie
- ▶ Potenciálne súvisiaci problém uvedený na GitHubu

2. Čím si vysvetľujete nekonzistentnosť v zarovnaní pôvodného a reverzného vlákna programom GraphAligner? Je k tomuto niečo v dostupnej literatúre, popr. kontaktovali ste tvorcov programu?



3. Ktoré z uvedených filtrácií hodnotíte ako najperspektívnejšie pre ďalšie použitie?

- ▶ Filtrácia homopolymérov

4. Aká bola dĺžka behu, bola by rovnaká schéma ako na Obr. 3.1 použiteľná aj pre väčšie genómy, napr. ľudský?

- ▶ Nie
- ▶ CPU servre
 - 47–142GB pamäte
 - 16–48 procesorov

