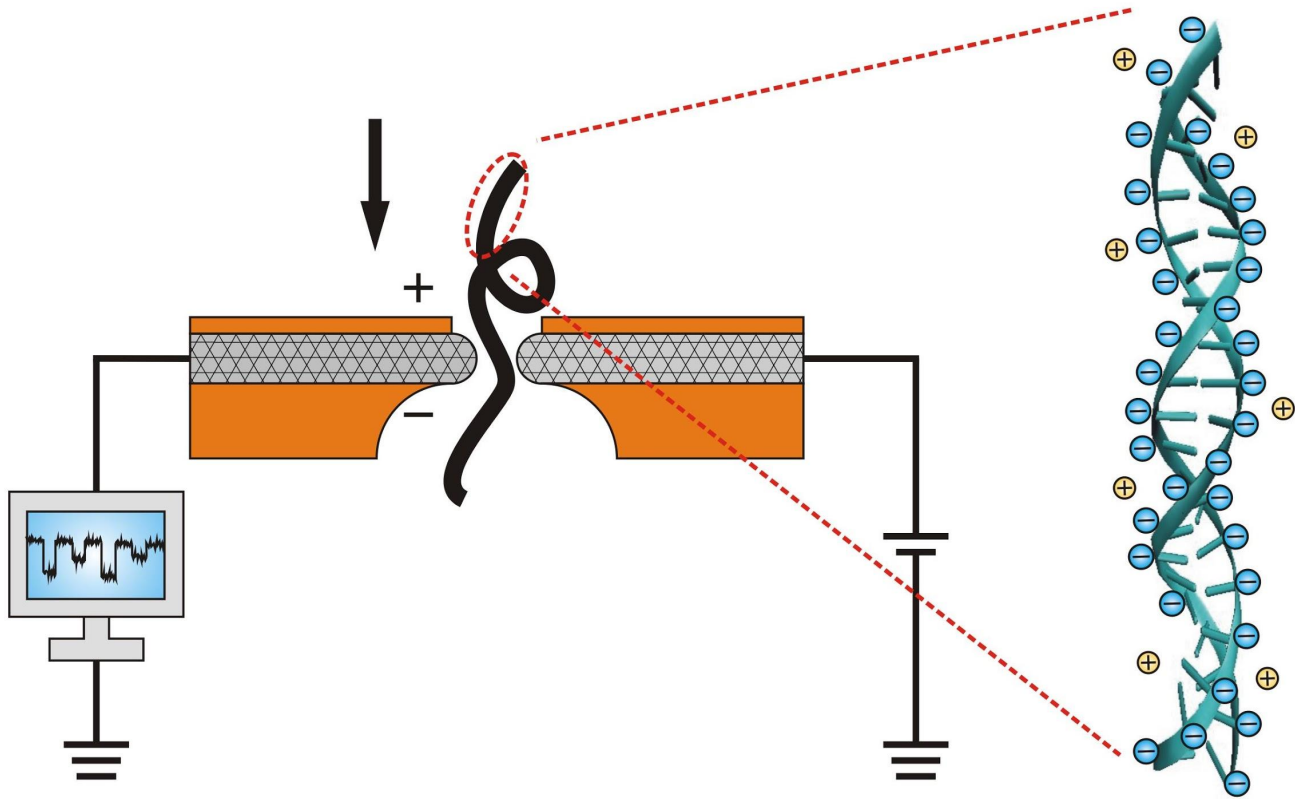




Encoding DNA-Sequenced Nanopore Data

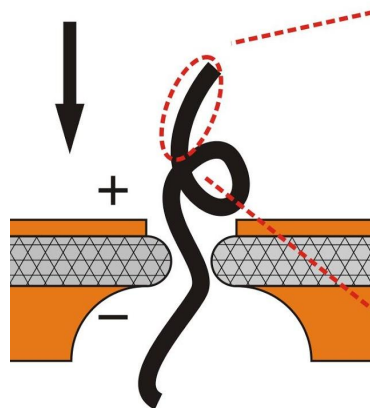
Presented by Sasha Jenner

Background

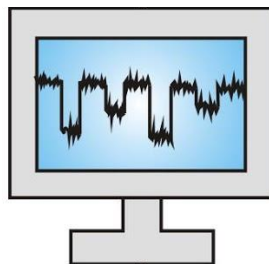


Nanopore sequencing of DNA

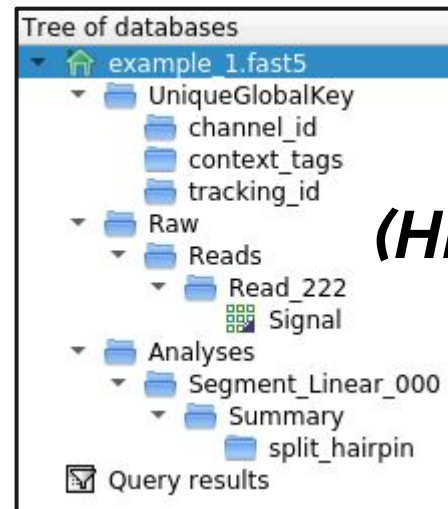
FAST5



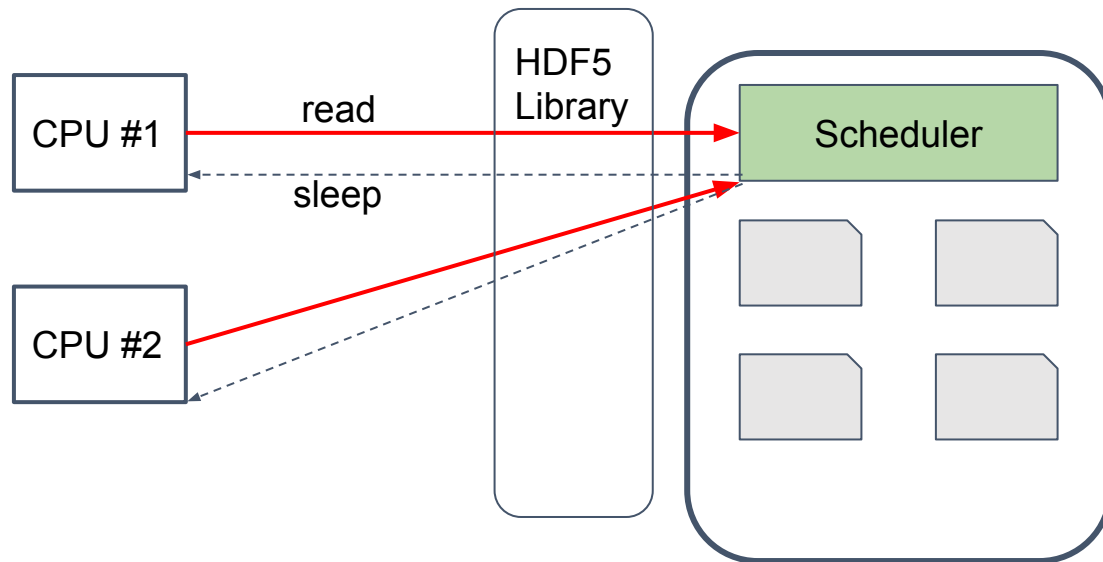
Read



Signal



(HDF5)



SLOW5 TSV File

```
#fileformat: slow5v1.0
#exp_start_time: 2020-01-01T00:00:00Z
#run_id: 855cdb4b26948
#flow_cell_id: FAH00000
```

Nanopore Metadata



#read_id	n_samples	digitisation	offset	range	sampling_rate	raw_signal
read-0	123456	8192	6	1467.6	100000	498,492,501,508,503,505,509,...
read-1	2000	8192	5	1467.6	4000	400,401,500,403,407,478,510,...
.
.
.
read-N	10000	8192	3	1467.6	4000	559,545,560,551,550,565,701,...

↑
Read ID

Read Metadata

↑
Signal

SLOW5 Index

	Location 	Read Size 
#read_id	file_offset	rec_length
read-0	67	500000
read-1	500001	1000000
•	•	•
•	•	•
•	•	•
read- N	364459005610	1580072

Aim

Smaller SLOW5 encodings

Methods

SLOW5 Encodings



In 1 Byte?

ASCII

0-9

VS

Binary

0-255

ASCII

example.blow5										9059, 1, 1		Top	
##file_format=slow5v0.1													
read_id		n_samples	digitisation	offset	range	sample_rate	raw_signal	num_bases	sequence	fast5_pa			
11b6cd19-3958-4264-a6f0-04aef956ebb		6028	8192.0	3.0	1467.6	4000.0	1373, 712, 738, 715, 716, 719, 728, 724, 727						
a649a4ae-c43d-492a-b6a1-a5b8b8076be4		59676	8192.0	23.0	1467.6	4000.0	1039, 588, 588, 593, 586, 574, 570, 585						
c3491225-815c-408b-abc6-ed864f545f4b		37454	8192.0	14.0	1467.6	4000.0	1099, 665, 655, 638, 622, 620, 627, 642						
7d717303c-726a-407f-8df6-59e98ef86e34		15665	8192.0	18.0	1467.6	4000.0	1235, 612, 611, 573, 581, 593, 572, 577						
9dc4d6c4-1dc0-49d0-aaa2-078408a749cf		52190	8192.0	6.0	1467.6	4000.0	1023, 606, 624, 546, 600, 581, 569, 564, 580						
52b95332-1cf5-4a6f-8bc4-88fbb1cb0c2c		44141	8192.0	3.0	1467.6	4000.0	1377, 956, 721, 682, 675, 695, 707, 700, 698						
8c395415-c8d4-4476-b77c-30c878bd8a1d		57421	8192.0	3.0	1467.6	4000.0	1173, 717, 722, 709, 714, 721, 712, 711, 718						
3fdd0b4a-2183-45ed-a817-c96e0b692df5		36568	8192.0	3.0	1467.6	4000.0	1370, 765, 680, 684, 696, 696, 684, 667, 617						
ca0779cd-f7a9-47ed-bd69-d50d61ce1c72		13002	8192.0	6.0	1467.6	4000.0	1257, 658, 527, 534, 533, 524, 531, 520, 527						
627a9fdf-1655-4b39-a413-c8f0dfb73dc6		45690	8192.0	5.0	1467.6	4000.0	1255, 773, 617, 574, 568, 555, 557, 554, 551						
-													
example.slow5										111, 12, 1		All	

All

<8>>^H^@^A^@^A^@^A^@^A^C%£1^N^O^L^FáUúFv^D^CGfby%N°<94>8%véöÖÁúö½<94><8a>6á2fGc·6íð{¿=<96>4^E<99>-5<93>±£iM<9c>>²V^Mu<
84>^N^EQ<8a>KÐ<84>Uj;à<89><90>«<9c>iZ^M<8d>iè¼Ä¯ōās<88>½<84>
<Vp^AQiá^Gô<91>w z^A^@^A^_<8b>^H^A^@^A^@^A^@^A^C5<9a>^G|<8f>x÷ÇóÜ«^a<8b>ö(-<9d>^ZU<9b>
~5j<8f><96><96>j<89>^R^S^f<83>§ö¨^A|^A^D^A<8c><90><84><90>h<8c>Ø3F<90>È^R
^R<91><90>aIMKñŋ^çf>ÂYü¼<9e><9c>qÍ¿<95>\\p̈jòĭ] ŊB<9a>®Böv-[yÜ:
µēōūō½<9e>g<8d>okôq÷ðt- [I£wi¾³É°üý ¥<8b>K<93>&Ó|v4wq^YQcçkGM;I6)&I<9c>6'M<96><89>3ûİ^F³A^Λ3ç¹*
İ5WUÖİ^Y<93>mr,{Ö\`KâyİÊ';İ<97> &<9d>ücM<90> O^KİL3Ëø<9b>@³<98>?^QwFro²<99>Ç}?
3<9d>{<93>İBy^_ËÜ_¼⁹<8d>}^F=ædv<88>Ycfś<9c>k<86><9b>¼ÆÈ^L<80><9a>lú;n<8a>Ü ^|jciz;Ü.³^A<9c>-²^İ^n<83>i,
İ~+avµ<9d>È<93>üv<81>|b×Üiv)oz[?ÜÈv³^S^O^ç·^Wà7İ¼g3äyŁÜk^NÑr<90>ũİT4®ÆSôműöTá¼^FÜWÜ-c:
<9b>U&^Jp<99>ÁP·ÜLBĞŲ|ıv¼q7z<99>^&y-<82>âİđĖčD^M-LkÔİ'dkc~5¼M²oëø^F# /3ĐØ <9f><93><9b>È¶ĂtäN_ăĚó f*
\\ydi:Öpī:÷<9d>Êđi^G^Eè<90>^4rH7×İ^C<93>o¶AE>^S^E5;M,
ÜÜēō<98>xPū<85>^-<9b><9a>°|<99>io^Z<9b>F|<8b>ũNt4_<9a>Şİ^Kç@s<8b>ý^AÇÜİEçuzÜÄv<90>İ^O®,6<9b>^CHf?
ÇYpm<9b>³Y'İ<97>ë-İpèg6z^Lf^Oc^AV|^ün<89><89>D<8f>İmsó^V:Ü+u|7^_u
Zòù<8a><99>b~@ªİ·,<91>-^Rh9<8d>ùòDÜ><8a>L^FòU^?<95>e7Ó^İ^İMP<,0Èİ÷ð=Ü<99>=O/
Q&^TÚ^İ^F<98><8f>İühş#SU<94>2^_<98><8f>İ<87>|<8c>ũÄ41<9f>s·ª^Ayéō<99>)ÇY^Uó<85>)È^[_p\$Ų÷>D^Kóİzt½^L^ÄöçY!
(JE^BQP<92>^C%ŞÑÈ+ßō<86><9f>ÜHN^Ĵ^_<81><8b>SiōÈi<^TpłZÄx·ú^E|µ<87>j7ó;őwáØ^¼|g¼Guő9¶İ<9d>L+ôP^OÜÈ@
İ^Sc]<8e>oBQ^Aēō6H+^GÆ3<90>ô^Uú^N^A^E^Upæ<91>90<8f>;<91>Q<9a>ũ^Ĵ[-E<99>?@<92>;øē^O^Fçİ
¶h<95>ý^F4=^omİÄVæÂN^RAScõ^AZ^W{ÇÜ^A^AjQB·L2xôB<8a>İðáo|¶+İ^G¶<88>|Yp¥xxz<8f>ÄËóys#ôWR~^HfêÄydwđpük<8a>Ú^Mô?
^PTô<81><96>İHà^O4~<LEUE<90>/ôÖ^H-V^NE^V^Máló^Ĵ[8~f·ñãoē<89>>z·İ]sßŠŲ^Ĵpİ<89>¶YĐ²^[_^N<83><pp~^P-\$<80>']ěášWøØÄZou
İ^UEbÄ9¶B.fñĴ[ĴÜ<8~ç^G^XÉō<9f>YŃon^AY^AOāİº^A^[[<^AQ?<8A<82>· Ab^?<99>¥à'İ<86>-
<9a>}½⁹<81>ôd<81>~Å^H-^Q_5<86>ô^?<9d><93>N|sÜ®Æ^WEÁÄ"%X^AGMÑ c^Ac^ZUH~ēňf<9e>pĐ^Miô@ÜGNō68đE7^S<90>ô4đ1^Nİ<8f>çG<9
1><96>'İ&uo<86>D¼B2Uđ)<95>;â+[×ō³đİwl5ms<8e>_Ø<8a>İ<95>m<^<89>UÄNg,
ō½^Kêg"ç<83>øv2Ç^İ^SÄT^7<91>éqoy^İ^IUĐç^Thí<83><94>J<83>©<92>İµñ<^M8Plİ<8d>9w<85><8e>Fü.^@¥-Đj^Sİ^Ĵb<93>M@
İB³<92>m'<9c>^Ä«<9c>^AAQ<ô<95>
geäg,¼<8c>çÇKô<8b>ê-c<90>X^W^H^v<90>ý?Zn<87>N<9a><82>ñĎĒÁÄ50AÔBĐ»^P^Ĵ^K<85>K?d.éu,
^YBo-y+^Gējİ·ypy<96>pñ£Nēō|<86>÷İ<99><8b>Hx:øø^@.ˆABÉÁ<93><8b>¹İ<p^Ĵă+^L^ZÝfpj^K"~<8a>Dò<9b>ÑæF<90>²^T[p<8f>tōNJ,
İÄāâ *<81>7ô<81>·XøIFRÜD<İĔä3(^Y<88>|^R,w^XJv+İyâ^A^O«İÉx±<9c>İGËYac!
2ēō<97>İü<83>Ü<89><Iä<8d>^L<90>{L<91>»^A4İôġâ±;oia^DÜP·ô+¶T ½?
pn8^G<9d>t't'ÇEuôÄU^M^Ĵ<8f>İ<82>^EADÄM^Ų<8e><80>J<89>¹; <91>ũ^Đp^IXÉY^D^L<8a>5oS;Ü<86>^F\$^O^A^E¶3xk^T^Ĵ[gF<9b>?
ÄÜD<90>ç^AK^H¼¹ōçβ^Wk^Ĵ^YA^_Ä<97>/\$ò^ĴU-wÄ

example.clow5

9305, 1, 1

[Top](#)

Algorithm

Input



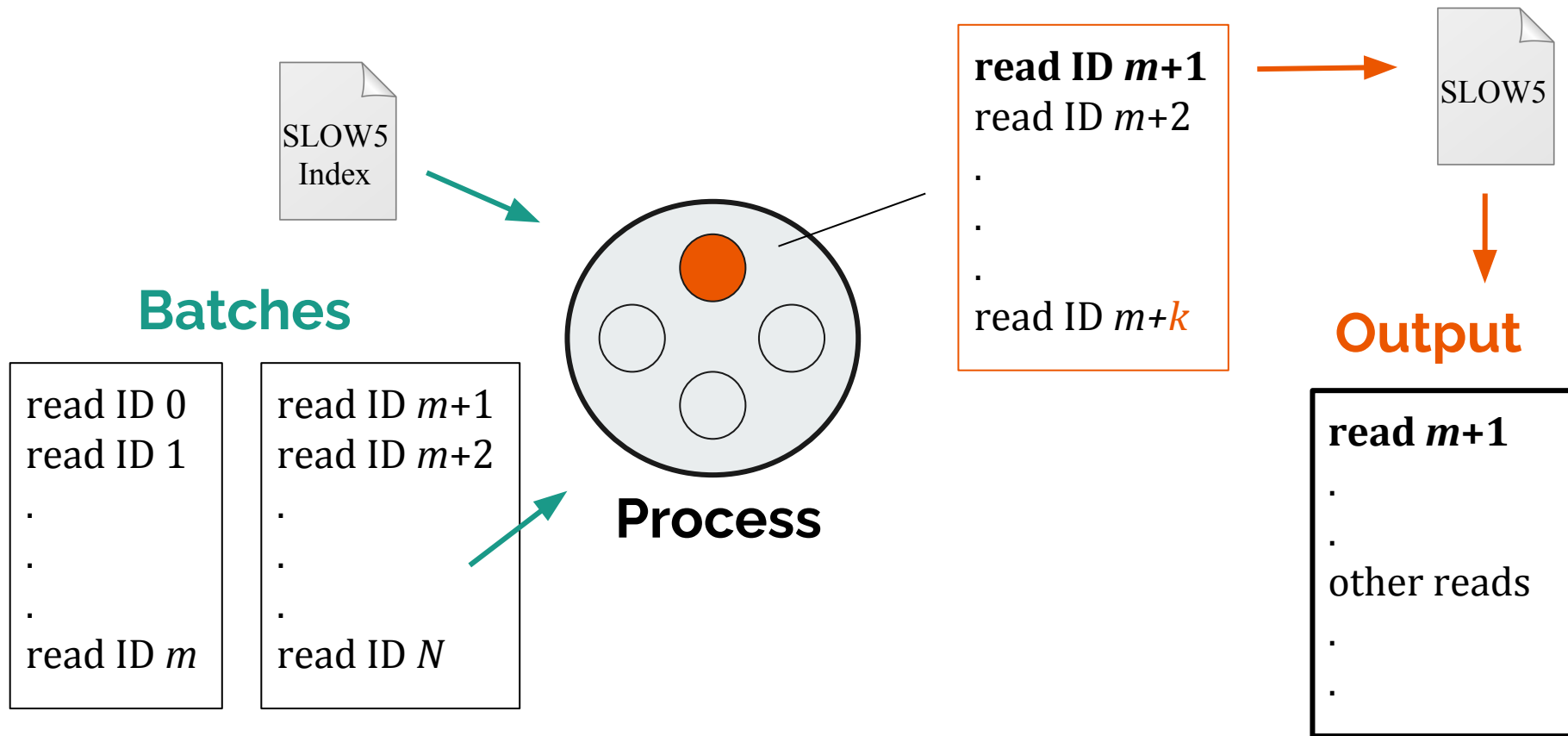
read ID 0
read ID 1
.
.
.
read ID N



Output

read 0
read 1
.
.
.
read N

Parallel Access

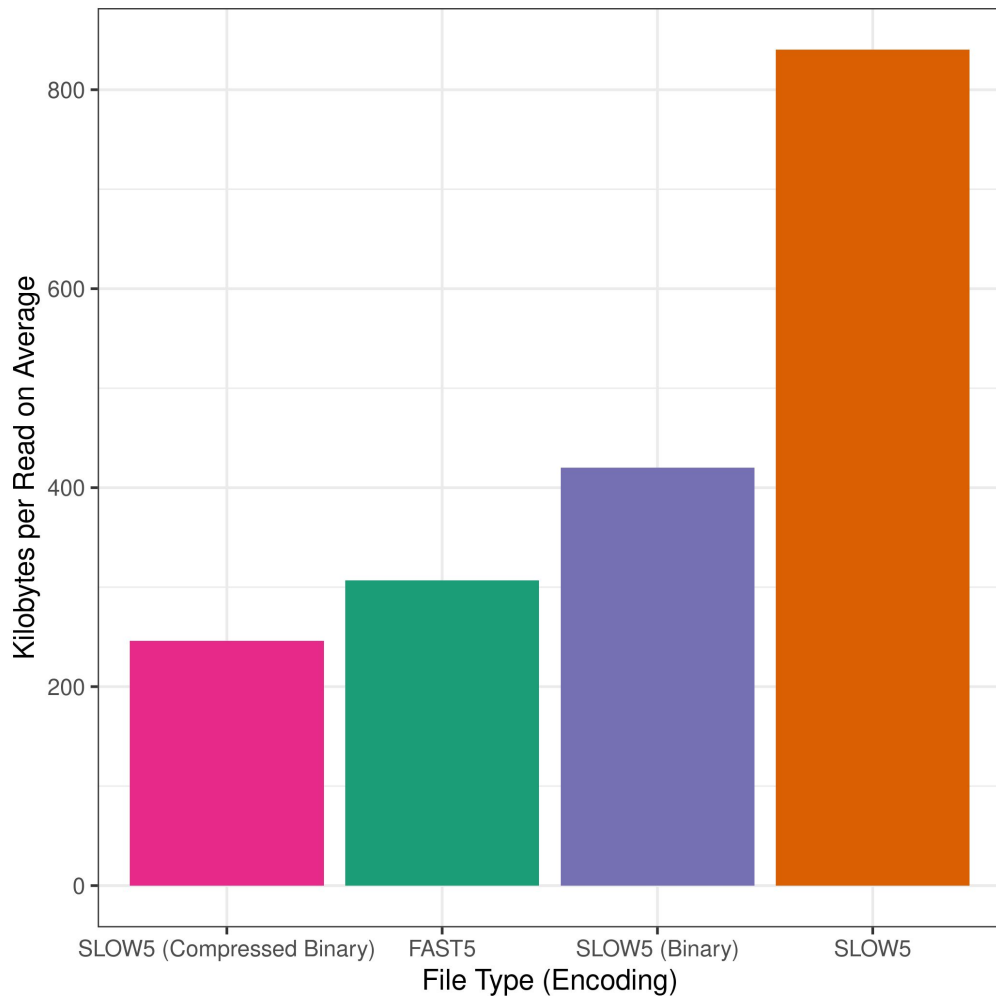


Results & Discussion



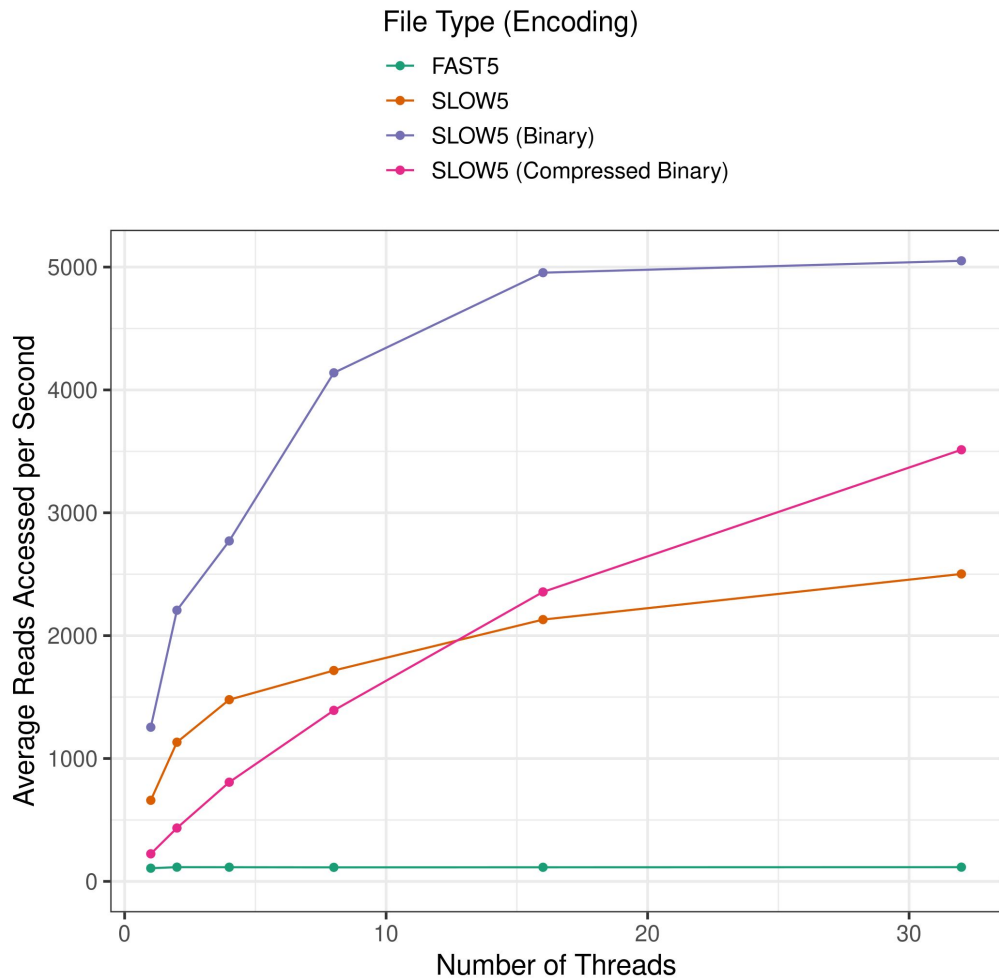
Benchmarking

1. File size
2. Access time



**Compressed
binary:**

**20% reduction
in FAST5 size**



E.g. 500 000 reads
FAST5: > 1h
SLOW5: < 5 min

Conclusion

- Compressed binary - most useful
- Binary encoding - fastest
- Both?