

Lab chat on my PhD

“Novel informatic sequence analysis approaches to dissect complex plant genomes”

April 2, 2015



What?

- ▶ Novel algorithms to analyse large-scale genomics data

What?

- ▶ Novel algorithms to analyse large-scale genomics data
- ▶ Our wish-list:
 - ▶ Reference & alignment free
 - ▶ Tolerates any sequencing platform
 - ▶ Works with Borevitz-style “wide and shallow” expts: e.g. 1000 samples at 1x

What?

- ▶ Novel algorithms to analyse large-scale genomics data
- ▶ Our wish-list:
 - ▶ Reference & alignment free
 - ▶ Tolerates any sequencing platform
 - ▶ Works with Borevitz-style “wide and shallow” expts: e.g. 1000 samples at 1x
- ▶ “...large data sets and high error rates combine to provide a ...challenge: it is now straightforward to generate data sets that cannot easily be analysed” — C T Brown *et al.* (2012)

What?

- ▶ Novel algorithms to analyse large-scale genomics data
- ▶ Our wish-list:
 - ▶ Reference & alignment free
 - ▶ Tolerates any sequencing platform
 - ▶ Works with Borevitz-style “wide and shallow” expts: e.g. 1000 samples at 1x
- ▶ “...large data sets and high error rates combine to provide a ... challenge: it is now straightforward to generate data sets that cannot easily be analysed” — C T Brown *et al.* (2012)
- ▶ e.g. Norman & I analysing ≈ 20 TB 3k rice genomes data. Our supercomputer allocation is 30 TB.

- ▶ *k*-mer analysis: analyse *k*-length words of sequence
 - ▶ Fast
 - ▶ Constant-memory (with `khmer`)
 - ▶ Scalable (linear time w/ number of samples)
 - ▶ Parallelisable (within & across nodes)

- ▶ *k*-mer analysis: analyse *k*-length words of sequence
 - ▶ Fast
 - ▶ Constant-memory (with `khmer`)
 - ▶ Scalable (linear time w/ number of samples)
 - ▶ Parallelisable (within & across nodes)
- ▶ Multi-layered “zooming” analysis
 - ▶ First pass basic clustering
 - ▶ Error correction
 - ▶ Population graph “alignment”
 - ▶ Variant calling

- ▶ *k*-mer analysis: analyse *k*-length words of sequence
 - ▶ Fast
 - ▶ Constant-memory (with `khmer`)
 - ▶ Scalable (linear time w/ number of samples)
 - ▶ Parallelisable (within & across nodes)
- ▶ Multi-layered “zooming” analysis
 - ▶ First pass basic clustering
 - ▶ Error correction
 - ▶ Population graph “alignment”
 - ▶ Variant calling
- ▶ In-silico experiment-driven development (new paradigm for me)

What have I been up to

- ▶ *k*-mer based clustering

What have I been up to

- ▶ *k*-mer based clustering
 - ▶ In collaboration w/ Sylvain Foret, Cheng-Soon Ong, Christfried Webers (NICTA)

What have I been up to

- ▶ *k*-mer based clustering
 - ▶ In collaboration w/ Sylvain Foret, Cheng-Soon Ong, Christfried Webers (NICTA)
 - ▶ Extending work in alignment-free sequence comparison (SF) and text/document clustering (C-SO, CW @ NICTA).

What have I been up to

- ▶ *k*-mer based clustering
 - ▶ In collaboration w/ Sylvain Foret, Cheng-Soon Ong, Christfried Webers (NICTA)
 - ▶ Extending work in alignment-free sequence comparison (SF) and text/document clustering (C-SO, CW @ NICTA).
 - ▶ Using Titus Brown's *khmer* (contributed a lot of code myself)

What have I been up to

- ▶ *k*-mer based clustering
 - ▶ In collaboration w/ Sylvain Foret, Cheng-Soon Ong, Christfried Webers (NICTA)
 - ▶ Extending work in alignment-free sequence comparison (SF) and text/document clustering (C-SO, CW @ NICTA).
 - ▶ Using Titus Brown's *khmer* (contributed a lot of code myself)
 - ▶ Initial results mixed: promising & and bit depressing

- ▶ *k*-mer based clustering
 - ▶ In collaboration w/ Sylvain Foret, Cheng-Soon Ong, Christfried Webers (NICTA)
 - ▶ Extending work in alignment-free sequence comparison (SF) and text/document clustering (C-SO, CW @ NICTA).
 - ▶ Using Titus Brown's *khmer* (contributed a lot of code myself)
 - ▶ Initial results mixed: promising & and bit depressing
 - ▶ New C++ implementation: Can compute 100x100 matrix in 10 mins on 16 CPUs
 - ▶ Hashing takes ≈ 10 mins, uses ≈ 1 GB/sample.

- ▶ *k*-mer based clustering
 - ▶ In collaboration w/ Sylvain Foret, Cheng-Soon Ong, Christfried Webers (NICTA)
 - ▶ Extending work in alignment-free sequence comparison (SF) and text/document clustering (C-SO, CW @ NICTA).
 - ▶ Using Titus Brown's *khmer* (contributed a lot of code myself)
 - ▶ Initial results mixed: promising & and bit depressing
 - ▶ New C++ implementation: Can compute 100x100 matrix in 10 mins on 16 CPUs
 - ▶ Hashing takes ≈ 10 mins, uses ≈ 1 GB/sample.
 - ▶ Norman, Sylvain, Cheng-Soon and Chris working on new metrics when SF & I are back