

# proovread manual

Thomas Hackl

## Contents

### 1 Installation

```
git clone --recursive https://github.com/BioInf-Wuerzburg/proovread
cd proovread/util/bwa
make
```

### 2 Dependencies

- Log::Log4perl
- blastn (NCBI Blast-2.2.24+ or later)

proovread is distributed with binaries of shrimp and blasr as well as bwa source code. If you want to employ other installed version of these mappers have a look at the *Advanced Configuration* section.

### 3 Usage

Test your installation by running proovread on the included sample data set.

```
proovread --sample --pre /tmp/pr-sample
```

Don't run proovread on entire SMRT cells directly, it will only blast your memory and take forever. Split your data in handy chunks of a few Mbp first:

```
SeqChunker -s 20M -o pb-%03d.fq pb-subreads.fq
```

Run proovread on one chunk first.

```
proovread -l pb-001.fq -s reads.fq [-u unitigs.fa] --pre pb-001
```

If things go smoothly, submit the rest.

```
xargs  
qsub  
sbatch
```

## 4 Input

long reads

short reads MiSeq/HiSeq

unitigs

nanopore

custom scenarios

## 5 Output

By default, proovread generates six files in the output folder:

<code>.trimmed.f[aq]</code>	high accuracy pacbio reads, trimmed for uncorrected/low quality regions
<code>.untrimmed.fq</code>	complete corrected pacbio reads including un-/ poorly corrected regions
<code>.ignored.tsv</code>	ids of reads and the reason for excluding them from correction
<code>.chim.tsv</code>	annotations of potential chimeric joints clipped during trimming
<code>.parameter.log</code>	the parameter set used for this run

If you are interested in mappings (SAM) and other intermediary files from iterations have a look at `-keep-temporary`.

## 6 Advanced Configuration

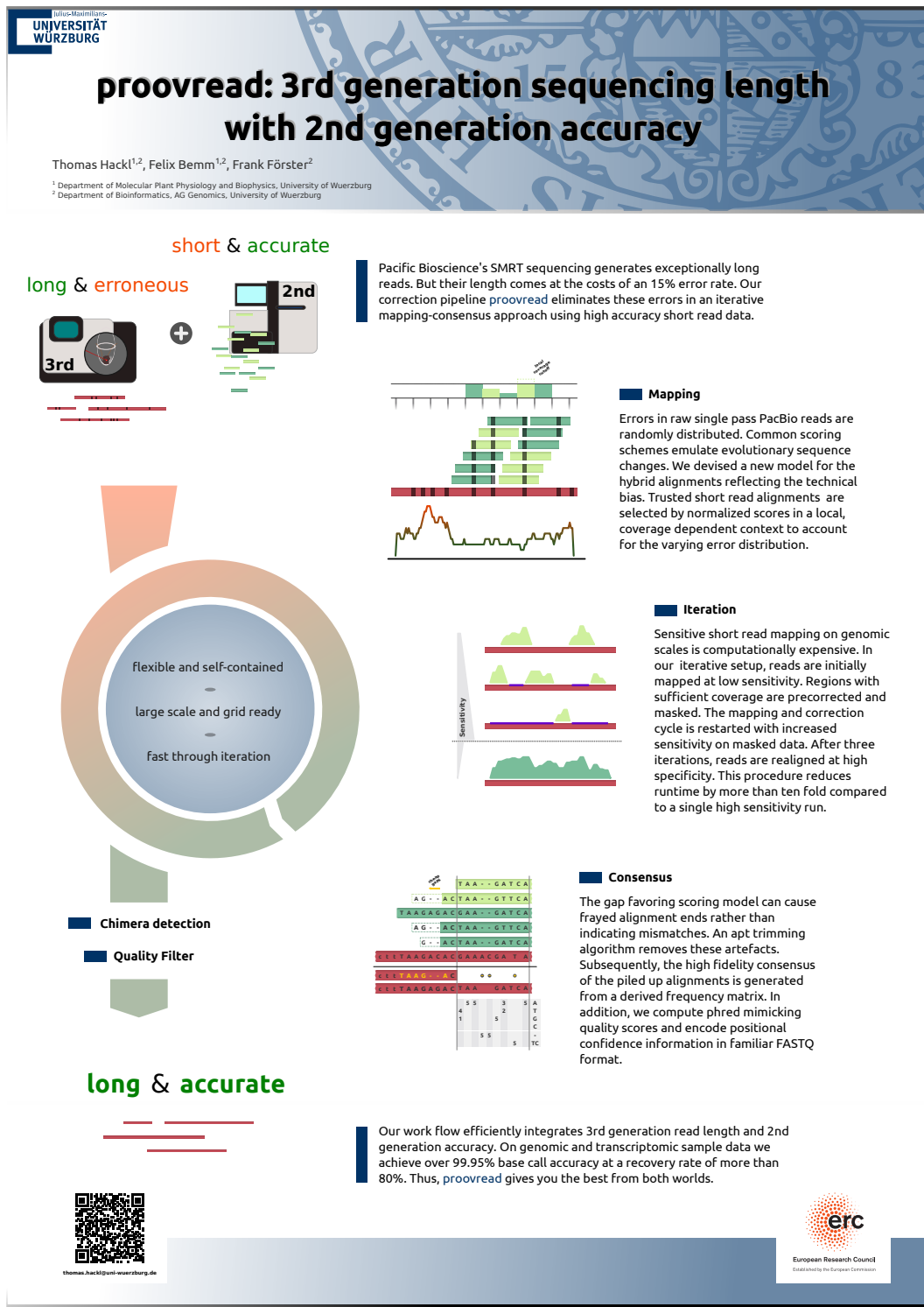
`-create-cfg`

## 7 Hardware and Parallelization

proovread has been designed with low memory node cluster architectures in mind. Peek memory is mainly controlled by the amount of input long reads. With chunks of less than 20 Mbp it easily runs on a 8 GB RAM machine.

Efficient parallelization in most cases is only possible for up to 4 or 8 threads in one instance.

## 8 Workflow



## 9 Citing proovread

proovread: large-scale high accuracy PacBio correction through iterative short read consensus. Hackl, T.; Hedrich, R.; Schultz, J.; Foerster, F. (2014).

<http://dx.doi.org/10.1093/bioinformatics/btu392?>

## 10 Contact

If you have any questions or encounter problems or potential bugs, don't hesitate to contact us. Either report issues on github (<https://github.com/BioInf-Wuerzburg/proovread/issues>) or write an email to:

- Thomas Hackl: [thomas.hackl@uni.wuerzburg.de](mailto:thomas.hackl@uni.wuerzburg.de)
- Frank Foerster: [frank.foerster@biozentrum.uni-wuerzburg.de](mailto:frank.foerster@biozentrum.uni-wuerzburg.de)