# *doepipeline*: a systematic approach to optimizing multi-level and multi-step data processing workflows

Daniel Svensson[1,$], Rickard Sjögren[1,2,$], David Sundell[3], Andreas Sjödin[3], Johan Trygg[1,2,*]

1. Department of Chemistry, Computational Life Science Cluster (CLiC), Umeå University, Umeå, Sweden

2. Corporate Research, Sartorius AG, Umeå, Sweden

3. Division of CBRN Security and Defence, FOI - Swedish Defence Research Agency, Umeå, Sweden

$ These authors contributed equally to this work

* Corresponding author

Johan Trygg: johan.trygg@umu.se

# 12 Abstract

13 **Background**

14 Selecting the proper parameter settings for bioinformatic software tools is challenging. Not

15 only will each parameter have an individual effect on the outcome, but there are also potential

16 interaction effects between parameters. Both of these effects may be difficult to predict. To

17 make the situation even more complex, multiple tools may be run in a sequential pipeline

18 where the final output depends on the parameter configuration for each tool in the pipeline.

19 Because of the complexity and difficulty of predicting outcomes, in practice parameters are

20 often left at default settings or set based on personal or peer experience obtained in a trial and

21 error fashion. To allow for the reliable and efficient selection of parameters for bioinformatic

22 pipelines, a systematic approach is needed.

23 **Results**

24 We present *doepipeline*, a novel approach to optimizing bioinformatic software parameters,

25 based on core concepts of the Design of Experiments methodology and recent advances in

26 subset designs. Optimal parameter settings are first approximated in a screening phase using a

27 subset design that efficiently spans the entire search space, then optimized in the subsequent

28 phase using response surface designs and OLS modeling. *doepipeline* was used to optimize

29 parameters in four use cases; 1) de-novo assembly, 2) scaffolding of a fragmented genome

30 assembly, 3) k-mer taxonomic classification of Oxford Nanopore Technologies MinION

31 reads, and 4) genetic variant calling. In all four cases, *doepipeline* found parameter settings

32 that produced a better outcome with respect to the characteristic measured when compared to

33 using default values. Our approach is implemented and available in the Python package

34 *doepipeline*.

35 **Conclusions**

36 Our proposed methodology provides a systematic and robust framework for optimizing

37 software parameter settings, in contrast to labor- and time-intensive manual parameter

38 tweaking. Implementation in *doepipeline* makes our methodology accessible and user-

39 friendly, and allows for automatic optimization of tools in a wide range of cases. The source

40 code of *doepipeline* is available at https://github.com/clicumu/doepipeline and it can be

41 installed through conda-forge.

42 Keywords

43 Design of Experiments, Optimization, Sequencing, Nanopore, MinION, Assembly,

44 Classification, Scaffolding

45

## Background

46

47 Bioinformatic software tools frequently offer a number of outcome-related parameters for the

48 user to set or change from their default values. These parameters may be different forms of

49 input filters, or alter the behavior of the running algorithm. Parameters may be either

50 quantitative or qualitative (multi-level) in nature. While it is advantageous to customize tools

51 to a specific situation, it is not always obvious what effect changing parameters will have on

52 the outcome. This may be due to lack of documentation, poor understanding of the algorithm,

53 or interaction effects between parameters that are difficult to foresee. Additionally, software

54 tools are commonly combined into pipelines, for example when calling genetic variants from

55 raw sequence reads [1,2]. Pipelining tools in this manner further increases the complexity of

56 selecting optimal parameter settings by increasing the numbers of both parameters and

57 potential interaction effects. The settings for a particular data processing pipeline may also

58 have to be tailored to the type of technology that was used to generate the data, for example

59 the different platforms available for DNA sequencing which yield different error profiles [3].

60 In general, the strategy for selecting parameter settings therefore typically consists of using

61 values derived from personal or peer experience and obtained in a trial-and-error fashion, or

62 simply retaining the default values. This kind of non-systematic selection of parameter

63 settings runs the risk of producing sub-optimal results.

64 The combined ranges of all possible parameter settings form a parameter space. To find the

65 optimal point in the parameter space, an exhaustive brute-force search, commonly called a

66 grid search, simply trying all possible combinations, is guaranteed to find the optimum. Since

67 the number of combinations increases exponentially, exhaustive searching quickly becomes

68 unfeasible as the number of parameters, and their ranges, grow. Instead, statistical Design of

69 Experiments (DoE) can be used to span and investigate the parameter space in an efficient

70 manner [4]. DoE aims to maximize information gain while minimizing the number of

71    experiments required [5]. This is done by introducing variation into the system under

72    investigation in a structured manner in order to explain how the parameters (*factors*)

73    influence the result (*response*). This variation is introduced according to statistical designs for

74    simultaneously varying the factor settings at a specific set of values (*levels*), and the system is

75    modeled using statistical methods, for example with Ordinary Least Squares (OLS)

76    regression [5–7]. The simplest type of statistical design is the full factorial design (FFD)

77    where all combinations of factor levels are investigated in an exhaustive manner, meaning

78    that they quickly become impracticable. To greatly reduce the number of experiments

79    required, fractional factorial designs (FrFD) are used to investigate a structured subset of the

80    FFD [6]. The problem is that FrFD are not trivial to use in situations where there are more

81    than two levels to investigate, and that there is no obvious way to combine qualitative and

82    quantitative variables. Recently, fractional factorial designs have been generalized into the so

83    called generalized subset design (GSD) [8]. GSDs are balanced and near-orthogonal multi-

84    level and multi-factor subset designs capable of mixing quantitative and qualitative factors,

85    allowing for the investigation of a large and diverse set of parameters in an efficient manner.

86    Compared to grid search, GSDs reduce the number of runs required to explore an equivalent

87    parameter space by an integer factor, also called the *reduction factor*.

88    Although DoE is primarily used in analytical chemistry, a DoE approach has previously been

89    applied by Eliasson et al to optimize software parameter settings in a liquid chromatography-

90    mass spectrometry (LC-MS) metabolomics data processing pipeline [9]. In essence, this

91    approach consists of sequentially updating a statistical design based on the predicted optimal

92    configuration of settings, until they converge at an optimum. We build upon the approach

93    proposed by Eliasson et al, and have developed a strategy for automated optimization of

94    software parameter settings. We extend Eliasson et al's approach with a screening phase

95    using the recently developed GSD to efficiently span a much larger parameter space. We also

96    make it possible to optimize multiple responses simultaneously. This extended approach may

97    be used both for optimization of individual tools and for multiple tools organized into a

98    pipeline. One crucial component is a well-defined objective function that you wish to

99    minimize or maximize, i.e. there must be some way to objectively determine how well the

100    pipeline is performing. Our strategy is software-agnostic and is implemented as a user-

101    friendly Python package - *doepipeline*.

102    In this article, we outline our DoE-based strategy for a systematic approach to optimizing

103    multi-level and multi-step data processing workflows, and exemplify the application of

104    *doepipeline* with four cases; 1) de-novo assembly of a bacterial genome, 2) scaffolding of

105    contiguous sequences (contigs) of a bacterial genome using 3rd generation sequencing

106    (nanopore) data, 3) k-mer taxonomic classification of long noisy sequence reads generated by

107    ONT MinION sequencing units, and 4) genetic variant calling in a human sample.

## Methods

109    We propose an approach for the optimization of software parameters, based on methods

110    derived from statistical design of experiments. Our approach, which has been implemented in

111    a python package (*doepipeline*), can be divided into two distinct phases:

112    1. Screening using a generalized subset design to find an approximate optimum. This

113        phase also serves to find the best choice of categorical variables.

114    2. Iterative optimization, starting from the best point found by screening, based on the

115        algorithm by Eliasson et al[9]. This phase optimizes only quantitative variables,

116        meaning that categorical variables are fixed at the best values found during phase 1.

117    The screening and optimization phases are schematically illustrated in Figure 1 and described

118    in more detail in the following subsections. Prior to screening and optimization, the user

119    specifies what parameters to use as factors in the designs, whether they are categorical or

6

120    numerical, and the permitted categories or value spans to be investigated. The user also

121    specifies what process outcomes to use as response, and whether it should be maximized,

122    minimized or reach a target value. In cases with several responses the user also needs to

123    specify low/high limits and the target for each response. The responses are then re-scaled

124    according to these limits and targets and combined into a single response using the geometric

125    mean according to Derringer & Suich desirability functions [10]. In brief, when there are

126    multiple responses each individual response is rescaled to be in the interval between 0 and 1,

127    and it is 0 when outside accepted limits and 1 when better than the target. The rescaled

128    responses are then combined into the overall desirability using the geometric mean.

129    **Figure 1 - Schematic visualization of *doepipeline* design space movement.** Example of

130    optimization of two factors (A and B) through both the screening and the optimization phase,

131    completed in 3 iterations. Each dot represents an executed pipeline with the parameters set by

132    factors A and B. Triangles represent executed pipelines using the optima of an Ordinary

133    Least Squares (OLS) model calculated in each optimization iteration. Red dots and triangles

134    represent the best configuration of factors found in each iteration. Dashed lines represent the

135    current high and low parameter settings in each iteration. *Screening phase*: a GSD using three

136    levels and a reduction factor of 2 is used to span the design space. The pipelines are executed

137    with the factor configurations suggested by the GSD and an approximate optimum is found

138    (red dot). *Optimization phase*: in iteration 2, an optimization design is created around the best

139    configuration found in the screening phase (black dots). In iteration 3, the design space is

140    moved in the direction of the configuration of factors that produced the best result (red

141    triangle) in iteration 2. *doepipeline* halts when the best response is produced by a

142    configuration of factors that lies close to the center point (red triangle in iteration 3).

7

**Screening for Approximate Optimum**

The purpose of the screening phase is to span the full search space to find regions with close

to optimal performance. Screening is performed by executing the specified pipeline using

combinations of factor configurations given by a GSD. Using a GSD effectively reduces the

number of experiments to run, while optimally spanning the search space (Fig. 1a). The

number of experiments required to investigate a given set of factors at a number of levels is

approximately an integer fraction of the total number of possible combinations, which

depends on the number of factors and their levels. A greater number of levels increases the

resolution of the space searched during screening but also exponentially increases the number

of runs required. We have found that five levels per numeric factor span large search spaces

with a high enough resolution to give satisfactory results, but it is possible to set the number

of levels individually for each factor in *doepipeline*. Similarly, we have found that the integer

fraction of the full design that the GSD should represent can be safely set at the number of

factors included in the design. However, this may also be controlled by the user by means of

the *reduction factor* setting in *doepipeline*.

The screening phase also serves the purpose of setting the category to use for each categorical

variable. For subsequent optimization, qualitative factors are fixed at the category of the best

factor configuration according to the screening. By fixing qualitative factors, only numeric

factors are investigated during the following optimization phase.

**Optimization of Numeric Factors**

After selecting the best factor configuration during screening, numerical factors are optimized

using response surface designs. The levels used in the screening design are here applied as

anchor points for the new optimization design. A response surface design, for instance a

central composite design, is constructed around the best configuration found. That is, the

8

167    configuration of factor levels found to produce the best result during the screening phase is

168    initially set as the center point in the new response surface design (Fig. 1b). If this

169    configuration lies at the edge of a factor's global design space (as defined by its min and max

170    allowed values), the factor's center point is shifted to the nearest screening level instead. This

171    is done in order to keep the design within the global design space. After having set the center

172    point for the new design, the high and low settings for each numeric factor are set to lie at the

173    midpoints between the nearest screening levels respectively above and below the chosen

174    center point, as indicated by the dashed line in Figure 1b. The span of each factor is then

175    defined as the difference between the high and low settings. As during screening, the

176    specified pipeline is executed using factor configurations given by the response surface

177    design.

178    During each optimization iteration, pipeline performance is approximated using OLS

179    regression [7]. By fitting a regression model the optimal configuration can be found by

180    optimizing the response predicted according to the model. The factors included in the OLS

181    model are selected either using a best subset approach or by using greedy forward selection;

182    the latter is preferred when more than four factors are included in the design. If the predictive

183    power (Q2) of the model is acceptable (Q2 > 0.5), the model is used to predict an optimal

184    parameter configuration. Each numeric factor's settings are then updated based on the best

185    result in a manner similar to the algorithm given by Eliasson et al [9]. For each factor, the

186    difference between the predicted best factor setting and the factor center point is calculated. If

187    this distance is greater than 25% of the span of the factor, the high and low settings of the

188    factor are updated in the direction of the best result. The default step length is 25% of the

189    span of the factor, i.e. the high and low settings are moved 25% of the step length (Fig. 1b,

190    iteration 3). We found that the algorithm did not always converge at this stage, but moved the

191    design space back and forth between iterations. To alleviate this problem, we implemented

9

192   design space shrinkage, which shrinks the design space span through multiplication by a so

193   called *shrinkage factor* (typical value is 0.9, corresponding to 10 % shrinkage) between

194   iterations, and found that it successfully improved convergence. If the proposed updated

195   factor settings lie outside the predefined design space limits, the design is instead placed at

196   the factor limits while keeping the same factor span. If the design has not moved between two

197   iterations, or the best response is not improved upon compared to the previous iteration, the

198   algorithm has converged and halts. If the optimization algorithm halts and responses have not

199   reached their minimally acceptable values, the screening results are re-evaluated and a new

200   optimization phase is run based on the results of the next best screening. At the end of the

201   optimization iterations the factor configuration that has produced the best result throughout

202   the iterations is chosen as the optimal configuration.

203   **Sequence data used in cases**

204   The *Francisella tularensis sp. holarctica* strain FSC200 [11]**,** and a genetic near neighbor

205   *Francisella hispaniensis* strain FSC454 were chosen as an example dataset in case 1 to 3 of

206   this study. The genome assembly of FSC200 is available as RefSeq assembly accesssion

207   GCF_000168775.2 and genome assembly of FSC454 as RefSeq assembly accession

208   GCF_001885235.1. Previously, generated Illumina HiSeq reads of FSC200 are available as

209   NCBI SRA run SRR518502. This latter dataset was subsampled down to an estimated

210   coverage of 100X (1.9M 100bp reads) for use in case 1, subsampling was performed with

211   *seqtk* [12](v. 1.2-r94, installed through bioconda [13]).

212   New sequencing libraries were prepared from DNA extractions of the two bacterial strains

213   using the SQK-LSK108 Ligation Sequencing Kit according to the manufacturer's

214   specifications and then sequenced using a FLO-MIN107 MinION flow cell (Oxford

215   Nanopore Technologies, UK). MinION sequence reads for FSC200 are available as NCBI

10

216    SRA run SRR9290761, and for FSC454 as NCBI SRA run SRR9290851. Subsampling down

217    to 50 000 from 132 259 MinION reads for FSC200 and 15 000 from 15757 MinION reads for

218    FSC454 was performed with a custom script, and the sequences were trimmed to a maximum

219    length of 3000 bp as well as being sorted by length to increase classification speed.

**Case 1: de-novo assembly of a bacterial genome**

221    In this example, we optimize the paired-end sequence assembler ABySS [14,15] (v. 2.0.2,

222    installed through bioconda [13]) to assemble the genome of an isolate of *Francisella*

223    *tularensis* ssp. *holarctica* (FSC200). ABySS has a total of 27 different parameters that can be

224    specified by the user. Some are directly related to the running time and memory usage of the

225    software (such as number of threads to use or bloom filter size), while others are related to

226    the quality and/or characteristics of the resulting assembly (such as the size of k-mer or the

227    minimum mean k-mer coverage of a unitig). For this example, we focused on the latter type

228    of parameter. Hence, all parameters chosen to be part of the optimization were deemed to

229    have a potential effect on the resulting assembly. The chosen parameters were: size of k-mer

230    ($k$) (KMER), minimum mean k-mer coverage of a unitig ($c$) (MIKC), minimum alignment

231    length of a read ($l$) (MIAL), and minimum number of pairs required for building contigs ($n$)

232    (MIPA).

233    For this optimization we set the investigated factor space so that the default value for each

234    factor was included within the span of each factors' min and max values (Table 1). Although

235    central to the ABySS algorithm, there is no default value for the k-mer size parameter. But

236    since the value of the k-mer size is bounded by the actual read length it was still possible to

237    define the GSD search space in a satisfactory way. For purposes of comparison, however, we

238    considered a k-mer size of 31 to be the default setting. In this example we ran the initial

239    screening with a reduction factor of 8, and used Central Composite Face-centered (CCF)

11

240    designs in the following optimization iterations. We used a shrinkage factor of 0.9 (-*s*), and

241    set the model selection method (-*m*) to greedy to speed up model selection. All other

242    *doepipeline* settings were kept at default values.

243    There are many metrics that can be used to evaluate the quality of a de-novo assembly, and

244    which specific ones to use depends on what the assembly is to be used for [16,17]. Example

245    metrics include the number of resulting contiguous sequences (contigs), the amount of total

246    sequence covered by the assembly, and the N50 value. The latter is the length of the contig

247    that, when the contigs are ordered by size, spans the midpoint of the total assembly. Hence,

248    the N50 value can be viewed as an assessment of the quality of the assembly in terms of

249    contiguity.

250    We used the total size of the assembly (tSeq), the number of contigs (nSeq), and the N50

251    value as responses. Since this optimization contained multiple responses, it was necessary to

252    set low/high acceptable limits for each response, as well as target values to reach. The low

253    and high limits for the responses were set with respect to the result obtained using the default

254    settings with the same input data, meaning that the worst acceptable results are the default

255    results. The target for the tSeq response was set to the reference genome size for FSC200

256    [11], while the targets for the nSeq and N50 responses were set to values that were

257    considered achievable (Table 2).

258    The data input to ABySS consisted of the subsampled Illumina HiSeq 2500 sequence data for

259    FSC200 (see Sequence data used in cases). Prior to calculating the values for the responses

260    we applied a length-based filter to the assembly using *Fastaq* [18] (v. 3.17.0), keeping only

261    those contigs more than 1000 bp in length. This filter was also applied when calculating the

262    response from the pipeline using the default parameter configuration. This is done because

263    the very short contigs are typically made up of short repetitive sequences, and removing them

12

264    simplifies the assembly graph and calculations on it. The software *seqstats* [19] was used to

265    calculate the response values from the filtered assembly.

266    **Case 2: scaffolding of a bacterial genome assembly using long reads**

267    Assembling a genome with short reads typically results in a fragmented assembly, consisting

268    of a number of contigs. The way these contigs are connected with each other - in terms of

269    ordering, distance, and direction - remains unknown. The reason for the fragmentation is that

270    certain stretches of genomes have low complexity and are therefore impossible to resolve

271    with short reads. One way of stitching together the contigs of an assembly is by using paired

272    reads with long insert sizes, or - as is increasingly common - using long reads from, for

273    example, the Nanopore or PacBio platforms. The long reads have an increased chance of

274    spanning the low-complexity regions, effectively anchoring both ends of a pair of contigs

275    together and thus resolving the gap. The process of connecting contigs together is referred to

276    as scaffolding, and the resulting sequences are known as scaffolds.

277    SSPACE-LongRead [20] (SSPACE) uses long reads, such as those produced by the PacBio

278    or Nanopore platforms, to scaffold an assembly. When running the software, the user can

279    manipulate a total of six parameters that relate to the resulting scaffolds. We investigated

280    whether manipulating some of the parameters would yield a better result than that achieved

281    by running SSPACE (v. 1-1) with default parameter settings. We chose to optimize the

282    minimum alignment length to allow a contig to be included for scaffolding ($a$) (ALEN), the

283    minimum gap between two contigs ($g$) (GLEN), the maximum link ratio between the two

284    best contig pairs ($r$) (RRAT), and the minimum identity of the alignment of the long reads to

285    the contig sequences ($i$) (IDEN). As response, we maximized the N50 value of the resulting

286    scaffolded assembly.

287    We set the investigated space for the factors so that the default value for each factor was

288    included within the span of each factor's min and max values (Table 3). For the optimization

289    phase following the screening phase we chose to use a CCF design for the experiments. The

290    reduction factor for the GSD was kept at the default value, i.e. the number of factors in the

291    investigation, which in this case was 4. The model selection method (*-m*) was set to greedy

292    and the shrinkage factor (*-s*) to 0.9. All other *doepipeline* settings were kept at default values.

293    The input assembly had been constructed with ABySS [15] (v. 2.0.2) (k=71) and subjected to

294    a contig length filter (>1000 bp). It consisted of 94 contigs between 1,685 and 87,479 bp in

295    length, had an N50 of 27,549 bp, and totaled 1,800,912 bp prior to scaffolding. The assembly

296    was constructed from the FSC200 Illumina HiSeq 2500 sequence data (see Sequence data

297    used in cases). We include the assembly at the *doepipeline* github repository. The read set

298    used for scaffolding consisted of 132,258 nanopore reads of between 163 and 108,214 bp in

299    length (N50=679 bp), totaling 104,374,862 bp. *Seqstats* [19] was used to calculate the

300    response from the scaffolded assembly.

**Case 3: k-mer classification**

302    K-mer classification is a method used to assign taxonomic labels to short DNA sequence

303    reads [21]. The method requires a precomputed database of k-mers generated from previously

304    known and assembled genomes, for example all complete genomes in the NCBI database.

305    When classifying a sample, the k-mer set of each read is calculated and compared with the

306    database of known k-mers. The read is then assigned to the most specific taxonomic class

307    within the database using the highest scoring k-mer root-to-leaf classification path following

308    the taxonomic hierarchy. This method is implemented in, for example, the software package

309    Kraken [22].

14

310    Kraken also uses a least common ancestor method, which re-classifies reads that are assigned

311    to multiple taxonomic sub-classes under a parent node. A read with non-unique leaf

312    assignment will then be assigned to the least common ancestor where there is little or no

313    assignment conflict instead. The k-mer classification method implemented in Kraken can be

314    applied to longer error-prone reads even though it is optimized for short accurate reads.

315    However, it will be less accurate due to the different (higher) error frequencies and will

316    therefore generate an increased rate of false positives.

317    In this study we used the software KrakenUniq [23] (v. 0.5.2). KrakenUniq builds upon the

318    Kraken engine but additionally records the number of unique k-mers as well as coverage for

319    each taxon. Three factors were used in the optimization: precision (PRES), minimum k-mer

320    hits (MH) and a filter (FILT). We chose to use a CCF design in the optimization phase of

321    *doepipeline*, the model selection method (*-m*) was set to greedy, and the shrinkage factor (*-s*)

322    to 0.9. All other *doepipeline* settings were kept at default values. The F1 score (Eq. 1), which

323    is the harmonic mean of precision and recall**,** was used as response.

324    $$F1 = 2 \cdot \frac{(precision \cdot recall)}{(precision + recall)}$$    **(Eq. 1)**

325    The input data were nanopore sequenced reads from two *Francisella* species, a target,

326    *Francisella tularensis holarctica* (FSC200) and one near neighbor, *Francisella hispaniensis*

327    (FSC454). The dataset was reduced to contain 50,000 *F. tularensis* reads and 15,000 (max) *F.*

328    *holarctica* reads of maximum length 3000 bp, to increase the speed of classification and

329    reduce potential bias (see Sequence data used in cases).

330    **Case 4: genetic variant calling**

331    A single genetic difference with respect to a reference genome is referred to as a genetic

332    variant, and the process of identifying these variants from sequence data is referred to as

15

333    variant calling. Calling the simplest form of genetic variant, single nucleotide variants (SNV),

334    from standard Illumina paired-end data is considered trivial nowadays, with F1 scores

335    reaching 0.98 [24]. Because of this, we opted to optimize calling of short insertions and

336    deletions (indels), which are slightly more complex and are harder to call correctly [24].

337    We used raw sequence data and high-confidence genetic (or "truth") variant calls from a

338    single well-studied individual, commonly known as NA12878. The raw sequence data

339    (2x100 bp, 50X depth), which form part of the Illumina Platinum Genomes (PG) [25], were

340    retrieved from the European Nucleotide Archive (ENA), study accession ERP001960 (run:

341    ERR194147). The truth callset was a "hybrid" dataset, meaning it was produced by

342    combining callsets obtained with different technologies and methodologies [25–27] as

343    described in Krusche et al [28]. The truth set was downloaded from the PG GitHub repository

344    [29].

345    The genome analysis toolkit (GATK) best practices workflow [1,2] was used as a guide for

346    this variant calling case. Raw data processing was carried out in accordance with GATK best

347    practices up to the point of having analysis ready reads, after which *doepipeline* was applied

348    to optimize the variant calling and filtering steps. First, PICARD (v. 2.18.1) [30] was used to

349    convert the sequence reads (FASTQ format) into unmapped BAM format (uBAM) and to

350    mark Illumina adapters. We then mapped the reads to the hg19 reference (part of the GATK

351    resource bundle) using BWA-MEM (v. 0.7.15 -r1140) [31,32] and marked duplicates using

352    PICARD. Finally, Base Quality Score Recalibration (BQSR) was carried out using GATK (v.

353    3.8-1-0) [33] to obtain analysis-ready reads.

354    This case aimed to optimize variant calling and variant filtering, the remaining steps in the

355    GATK best practices after obtaining analysis-ready reads. The calling was carried out using

356    HaplotypeCaller, and the filtering was carried out using VariantFiltration, both tools within

357   GATK. HaplotypeCaller offers around 20 adjustable parameters while the VariantFiltration

358   tool expects custom-specified cutoffs for annotations in the variant call format (VCF) file.

359   GATK suggests four annotations by which to filter indels. In order to include a meaningful

360   number of parameters at each step we chose to optimize the two steps sequentially.

361   Performing sequential optimization allowed us to investigate 4 parameters for each step, 8 in

362   total. We first optimized the calling step while keeping the parameters in the filtering step at

363   their default settings. We then optimized the parameters for the filtering step using the output

364   from the highest scoring experiment in the first step. For the calling step we chose to

365   optimize the global assumed mismapping rate for reads (*globalMAPQ*) (GMQ), the minimum

366   base quality for calling (*mbq*) (MBQ), the minimum reads per alignment start

367   (*minReadsPerAlignStart*) (RAS), and the minimum confidence threshold for calling

368   (*stand_call_conf*) (SCC). For the filtering step we chose to optimize the quality by depth

369   (*QD*) (QD), the read position rank sum test (*ReadPosRankSum*) (RPRS), the Fisher test for

370   strand bias (*FS*) (FS), and the strand odds ratio (*SOR*) (SOR). To further reduce the size of

371   the optimization, we chose to optimize only against variants on chromosome 1. However, we

372   screened for any overfitting of the parameters by executing the variant calling and filtering

373   pipeline across all autosomes and chromosome X with the optimized parameters.

374   The following settings were used for both optimizations. We set the space investigated for the

375   factors so that the default value for each factor was included within the span of each factor's

376   min and max values (Table 5). The design of choice for the screening phase was the CCF

377   design. The reduction factor for the GSD was increased to 8, reducing the number of

378   experiments. The model selection method (*-m*) was set to greedy and the shrinkage factor (*-s*)

379   to 0.9. All other *doepipeline* settings were kept at default values.

380   Performance metrics and tools to assess the accuracy of variant calling in a standardized

381   manner are crucial, and the benchmarking team of the Global Alliance for Genomics and

17

382    Health (GA4GH) have made significant progress with respect to this [28]. The GA4GH

383    benchmarking team has developed a benchmarking tool, hap.py [34], that can compare a

384    high-confidence (or "truth") variant callset with a user-made single-sample callset, also

385    known as the query callset, and output performance metrics. For a certain set of confident

386    regions (specified by BED file), concordant variants in the two callsets should be considered

387    true positives (TP), while discordant variants should be considered either false positives (FP)

388    or false negatives (FN) depending on which callset they appear in. Hap.py also outputs the F1

389    score (see case 3 methods) for variants passing the VCF filters, which was used as the

390    response in this case.

391    **Grid search comparison**

392    We compared the results from *doepipeline* to those from grid search, which is a common

393    methodology for optimizing parameters. Grid search is done by evaluating the parameter

394    performance for all possible combinations of parameter settings, the so-called parameter grid.

395    For the comparison to be relevant, we performed the grid search at the same resolution as the

396    GSD screening step in each case. In other words, we tested all possible combinations of the

397    factor setting levels (typically 5 levels per factor).

398    *doepipeline*

399    **Implementation**

400    *doepipeline* is fully implemented in the Python programming language and source code is

401    available for download at github (https://github.com/clicumu/doepipeline) and installable

402    with conda-forge [35,36] and through PyPi. Generation of statistical designs is carried out

403    through the python package PyDOE2 [37], in which the GSD has been implemented.

**Usage**

404

405 Configuration of the optimization is done in a structured YAML file with sections for the

406 experimental design and for the pipeline steps (commands) to run. The design section

407 includes the names of the factors investigated and their min/max values (design space), the

408 responses and their goals (minimize/maximize), and the type of design to use in the

409 optimization phase. The pipeline section is where each individual pipeline step is specified.

410 In each iteration, *doepipeline* takes the pipeline steps as configured and substitutes the

411 parameters under investigation with the values given by the statistical design. A batch script

412 is created for each pipeline step, with any parameter values substituted, and the execution of

413 it is controlled by *doepipeline*. Pipeline steps are executed either in parallel mode, where all

414 experiments are run at the same time, or in sequential mode where each pipeline with all of

415 the steps is executed in sequence. For reference, we provide example YAML files at the

416 github repository.

417 Today, scientific data processing can include vast amounts of data and/or require substantial

418 computing power. In such cases, data processing is commonly performed on compute clusters

419 that typically use some queueing system in order to handle all user requests for resources. An

420 example of such a queueing system is the Slurm Workload Manager [38] (Slurm). To

421 accommodate users of compute clusters, we have implemented Slurm support for

422 *doepipeline*. If using Slurm, specify the Slurm options in the YAML file as you would when

423 running a regular Slurm job. The Slurm options are transferred by *doepipeline* to the batch

424 script which is then submitted to Slurm using *sbatch*.

425 After optimization, the parameter values suggested by *doepipeline* are saved in the working

426 directory for the optimization. Additionally, there is a rich log file that can be investigated to

427 follow the workflow.

19

# Results

**Case 1: de-novo assembly of a bacterial genome**

The goal of de-novo assembly is to combine raw sequence reads into a representation of an organism's genome, i.e. to obtain as contiguous a genomic sequence as possible. Due to the characteristics of the genome sequence itself, in combination with short reads, this process can be difficult. For example, sequence reads from less complex segments of the genome will map to more than one position, causing ambiguities that are not possible to resolve, and this in turn leads to fragmentation of the assembly.

One popular sequence assembler is ABySS [15], which provides 27 different user-controlled parameters. We set up an example for optimization of de-novo assembly software parameters using ABySS (see Methods section). *doepipeline* ran for two iterations before halting. Thus, the best response was obtained in the first iteration, in the GSD screening phase. The experimental sheet and corresponding response values from the GSD screening and iteration 2 are included as Additional file 1. Using the optimized parameter settings (Table 1), we obtained a 1.6% and 13.1% increase in the investigated responses tSeq and N50, and a 2.2% reduction in nSeq as compared to when abyss-pe was run with default settings (Table 2). Optimizing the parameters using the grid search option required 625 experiments to be run, and it resulted in the same combination of parameter settings as when using *doepipeline* (see Additional file 2 for grid search result). By comparison, doepipeline required 97 experiments to be run.

**Table 1 - Factors in the de-novo assembly case.** The four factors investigated in the de-novo assembly case are described below. The letter in parenthesis following the parameter name is the parameter used in the abyss-pe command line interface. Min and max values define the design space. *: There is no default value explicitly specified by the ABySS

20

452    documentation. However here we used a k-mer size of 31 for comparison purposes. **: This

453    refers to the square root of the median k-mer coverage, which is affected by the sequencing

454    depth and choice of k-mer size. The optimized values are the combination of factor values

455    that produced the best outcome, as found by *doepipeline*.

| Parameter | Abbr. | Type | Min | Max | Default | Optimized |
|---|---|---|---|---|---|---|
| Size of k-mer ($k$) | KMER | Ordinal | 20 | 90 | 31* | 38 |
| Minimum mean k-mer coverage of a unitig ($c$) | MIKC | Quantitative | 2 | 15 | sqrt(median)** | 8.5 |
| Minimum alignment length of a read ($l$) | MIAL | Ordinal | 20 | 60 | 40 | 30 |
| Minimum number of pairs required for building contigs ($n$) | MIPA | Ordinal | 5 | 15 | 10 | 15 |

456

457    **Table 2 - Responses in the de-novo assembly case.** The three responses that were measured

458    in the de-novo assembly case are described below. *: Responses that have the criterion

459    maximize have a low limit, and those with the criterion minimize have a high limit. **:

460    Default values are based on using a k-mer size of 31 and leaving all other parameters

461    unchanged.

| Response | Abbr. | Criterion | Low/high limit* | Target | Default** | Optimized |
|---|---|---|---|---|---|---|
| Total sequence in assembly (bp) | tSeq | Maximize | 1,830,000 | 1,894,157 | 1,835,427 | 1,864,165 |
| Number of contigs in assembly | nSeq | Minimize | 95 | 85 | 91 | 89 |
| N50 | N50 | Maximize | 28,000 | 35,000 | 28,149 | 31,847 |

462

21

463 **Case 2: scaffolding of a bacterial genome assembly using long reads**

464 Scaffolding is the process of connecting together contigs obtained from an assembly step. In

465 this example we aimed to optimize parameters for the scaffolding software package

466 SSPACE-LongRead [20], which relies on long reads to span the low-complexity regions that

467 are typically found between the contigs of an assembly. *doepipeline* ran for three iterations

468 before halting, obtaining the best result in the second iteration. The response values and

469 parameter settings investigated in each iteration are included in Additional file 3. The

470 response (N50) value obtained when using the default parameter settings was 1,141,889 bp.

471 Using the optimized parameter settings (Table 3) resulted in a 66.9% increase in the response

472 (1,905,883 bp). Optimizing the parameters using the grid search option required 625

473 experiments to be run, compared to 211 experiments using *doepipeline*, and it resulted in a

474 best N50 value of 1,868,309, which is slightly lower than the result obtained using

475 *doepipeline* (see Additional file 4 for grid search result).

476 **Table 3 - Factors in the scaffolding case.** The four factors investigated in the scaffolding

477 case are described below. The letter in parenthesis following the parameter name is the

478 parameter used in the SSPACE command line interface. Min and max values define the

479 design space. The optimized values are those that in combination produced the best outcome,

480 as found by *doepipeline*.

| Parameter | Abbr. | Type | Min | Max | Default | Optimized |
|-----------|-------|------|-----|-----|---------|-----------|
| Minimum alignment length to allow a contig to be included for scaffolding (*a*) | ALEN | Ordinal | 0 | 5000 | 0 | 0 |
| Minimum gap between two contigs (*g*) | GLEN | Ordinal | -3000 | 3000 | -200 | -750 |
| Maximum link ratio between two best contig pairs (*r*) | RRAT | Quantitative | 0.1 | 0.7 | 0.3 | 0.325 |

| Minimum identity of the alignment of the long reads to the contig sequences ($i$) | IDEN | Ordinal | 30 | 90 | 70 | 82 |
|---|---|---|---|---|---|---|

**Case 3: k-mer classification**

K-mer classification is used to gather information about the species content of a metagenomic sample. It is possible to visualize the general distribution of species through the reads classified or to identify the presence/absence of reads classified to specific targets. By using third generation sequencing techniques, such as Oxford Nanopore, it is possible to classify reads from an unknown sample in real time. But due to the long error-prone reads produced by third-generation sequencing machines, there is a greater risk of misclassification. At the genus level this is not usually a problem. But when it comes to discriminating between pathogenic and non-pathogenic species, misclassification may become problematic; in particular false positive signals of pathogenic species may be obtained. We investigated the KrakenUniq [23] (v. 0.5.2) algorithm and used *doepipeline* to find optimized settings for long error-prone reads in order to increase the ratio of true positives to false positives using the F1-score as response. KrakenUniq also has a filter that may reduce the number of false positive reads. The filter will adjust each assigned read up the tree until the desired threshold is met, where the threshold is the number of assigned k-mers divided by the number of unique k-mers in that category [23].

Optimization ran for three iterations before halting and the best results were found during the second iteration. The experimental sheet and corresponding response values from the GSD screening and optimization iterations are included in Additional file 5. Using the optimized parameter settings (Table 4), we were able to increase the F1 score by 0.065% from 0.993690 to 0.994341, compared to when running KrakenUniq with default settings.

23

502  Optimizing the parameters using the grid search option required 125 experiments to be run,

503  compared to 76 experiments using *doepipeline*. The grid search resulted in a best F1 score of

504  0.994169 which is slightly lower than the result obtained using *doepipeline* (see Additional

505  file 6 for grid search result).

506  **Table 4 - Factors in the k-mer classification case.** The three factors investigated in the k-

507  mer case are described below. Min and max values define the design space. The optimized

508  values are those that in combination produced the best outcome, as found by *doepipeline*. *:

509  The KrakenUniq documentation to our knowledge does not state what the default value is.

| Parameter | Abbr. | Type | Min | Max | Default | Optimized |
|-----------|-------|------|-----|-----|---------|-----------|
| Minimum k-mer hits | MH | Ordinal | 1 | 200 | * | 14 |
| Standard deviation of the relative errors of the estimate | PRES | Ordinal | 10 | 18 | 12 | 17 |
| Minimum tax-ID score threshold | FILT | Quantitative | 0 | 0.05 | 0 | 0 |

510

511  **Case 4: genetic variant calling**

512  Variant calling is the process of determining genetic variants (or mutations) from genetic

513  sequence data. In this case we aimed to find optimized parameters for a widely used variant

514  calling framework, the genome analysis toolkit (GATK). Specifically, we sequentially

515  optimized two of the steps carried out by GATK: variant calling and variant filtering (see

516  methods).

517  In the optimization for the first step (variant calling), *doepipeline* ran for three iterations

518  before halting, obtaining the best result in the second iteration (F1=0.9707). In the

519  optimization for the second step (variant filtration), *doepipeline* ran for four iterations before

520  halting. The best result (F1=0.9716) was obtained in the fourth iteration when the optimum

24

521    predicted by the model was validated. This optimum was too far from the design space edges

522    for *doepipeline* to move the design space and initiate another iteration, and thus it halted

523    execution. The response values and parameter settings investigated in each iteration are

524    included in Additional file 7 (variant calling step) and 8 (variant filtering step). The included

525    parameters and their default and optimized settings are listed in Table 5.

526    As the optimization was performed only on chromosome 1, we wanted to see how well the

527    optimized parameter settings carried over into a variant calling and filtering pipeline applied

528    across all autosomes and chromosome X. This analysis resulted in an F1 score of 0.9713,

529    while using the default settings resulted in an F1 score of 0.9702.

530    Optimizing the parameters using the grid search option resulted in a best F1 score of 0.9715,

531    which is marginally lower than the results obtained using *doepipeline*. Five experiments in

532    the first step of the grid search optimization (calling) resulted in the same highest F1 score

533    (see Additional file 9). We therefore ran five parallel instances of the second step of

534    optimization (filtering) using the different VCF files from the five best experiments in the

535    first step. This inflated the number of required experiments from the expected 1250 to 3750

536    experiments in total, compared to 280 experiments with *doepipeline*. The five parallel

537    optimizations of step two all yielded the same set of 12 combinations of settings producing an

538    equally high F1 score (0.9715) (see Additional file 10). Validation across all autosomes and

539    chromosome X using all 60 combinations of parameter settings (5 times 12) yielded a best F1

540    score of 0.9712, again marginally lower than for doepipeline.

541    **Table 5 - Factors in the variant calling case.** The factors investigated in the variant calling

542    case are described below. The optimization was carried out sequentially for two main steps,

543    variant calling and variant filtering, and which step each factor belongs to is indicated. For

544    the variant calling step, the factor's corresponding command line flag is given in parentheses

545    after the parameter name. For the variant filtering step, the corresponding information tag

25

546  annotated in the VCF file is indicated in parentheses. The min and max values define the

547  design space. The default values for all factors are also indicated; for the calling step they are

548  the build-in default values of the HaplotypeCaller tool, while for the filtering step the default

549  values are those recommended by the GATK team. The optimized values are those that in

550  combination produced the best outcome, as found by *doepipeline*.

| Step | Parameter | Abbr. | Type | Min | Max | Default | Optimized |
|---|---|---|---|---|---|---|---|
| Variant calling | Global assumed mismapping rate for reads (*globalMAPQ*) | GMQ | Ordinal | 20 | 55 | 45 | 46 |
| | Minimum base quality for calling (*mbq*) | MBQ | Ordinal | 5 | 25 | 10 | 10 |
| | Minimum reads per alignment start (*minReadsPerAlignment*) | RAS | Ordinal | 5 | 25 | 10 | 20 |
| | Minimum confidence threshold for calling (*stand_call_conf*) | SCC | Quantitative | 5 | 25 | 10 | 5 |
| Variant filtering | Quality by depth (*QD*) | QD | Quantitative | 0 | 10 | 2 | 0.41 |
| | Read position rank sum test (*ReadPosRankSum*) | RPRS | Quantitative | -40 | 0 | -20 | -37.5 |
| | Fisher test for strand bias (*FS*) | FS | Quantitative | 0 | 250 | 200 | 62.5 |
| | Strand odds ratio (*SOR*) | SOR | Quantitative | 0 | 20 | 10 | 8.16 |

551

# Discussion

553  Selecting parameter settings for a data processing pipeline is complex, since the influence of

554  the parameters on the end result is not always obvious. While the value of personal and peer

555  experience should not be underestimated, our approach provides a systematic way of

26

556  determining optimal settings. Specialized tools to optimize particular bioinformatic software

557  tools have been proposed previously. For example, VelvetOptimizer [39] can be used to

558  optimize the k-mer and coverage cutoff parameters of the Velvet assembler [40] and

559  KmerGenie can be used to make an informed decision on the choice of k-mer in de Bruijn

560  based assemblers [41]. However, a generalized, software-agnostic optimization approach is

561  preferable, especially when several tools are used together in a pipeline.

562  Here we present such a generalized strategy for automated sequential optimization of

563  software parameters, employing core concepts of DoE methodology. We have implemented

564  our strategy in a user-friendly python package, *doepipeline*. The optimization strategy and the

565  use of *doepipeline* was exemplified in four bioinformatics use cases; de-novo assembly of a

566  bacterial genome using Illumina reads, scaffolding a bacterial genome assembly using

567  nanopore reads, k-mer classification of metagenomic third generation sequencing data, and

568  genetic variant calling. In all four cases, we saw an improvement in the measured response

569  variables as compared to when using the default parameter settings. The improvement of the

570  measured responses in our examples ranged between 0.065% and 66.9%. We compared the

571  results from *doepipeline* to results from standard grid searches, and *doepipeline* achieved

572  equally good or better results using significantly lower numbers of evaluations/experiments.

573  Grid search is typically limited to running a single optimization phase evaluating all points in

574  the parameter grid with no further refinement. This is in contrast to *doepipeline*, which is

575  adaptive and refines the parameter settings based on the best results from the previous phase,

576  allowing it to find better performing parameter settings than grid search.

577  One of the advantages of our proposed strategy is the use of a GSD in a screening phase prior

578  to the optimization phase. Compared to Eliasson et al[9], we are able to screen a much larger

579  design space efficiently prior to optimization using the GSD-based approach. In order for the

580  optimization phase to converge in a feasible number of iterations, the design space should be

581   restricted in some way. Deciding the range of each of the factors without guidance risks

582   creating too narrow or wide a design space. Instead, the screening allows the user to set up a

583   relaxed (wide) design space in which to investigate and to approximate the optimal factor

584   combination. The approximation represents a substantiated initial center point around which

585   to set up a narrower optimization design. The screening phase will also identify promising

586   values for any qualitative factors and fix them before optimization. Thus, the GSD screening

587   phase can be viewed as a systematic approach to restricting the design space for the

588   subsequent optimization phase. Similar results can be achieved using stochastic optimization

589   methods such as random search [42], commonly applied within the machine learning

590   community. Random search can effectively reduce the number of runs required, but the final

591   results are probabilistic and may not be optimal, depending on each particular random draw.

592   By using structured space-filling designs, *doepipeline* deliberately spans more of the search

593   space rather than relying on randomness. We note that the multi-phase workflow of

594   *doepipeline* has conceptual similarities to Bayesian hyperparameter optimization [43], in

595   refining the parameter choice based on promising parameter regions from earlier iterations.

596   However, *doepipeline* uses statistical designs that are guaranteed to fill the parameter space

597   and structured refinement around promising points rather than randomly sampling promising

598   regions with higher probability.

599   The fraction of the full design that a GSD represents can be controlled with the reduction

600   factor parameter in *doepipeline*. We ran the optimization of ABySS (case 1) with a GSD

601   reduction factor of 8, but another optimization of ABySS where a reduction factor of 10 was

602   used produced the same response values (data not shown) in fewer experimental runs (45 as

603   opposed to 70). In addition, there was a degree of overlap among the response values in the

604   GSD iterations (Additional files 1, 3, and 5). Overall, this could indicate that it is meaningful

605 to try running the GSD with a higher reduction factor than the recommended default, and/or

606 reducing the number of levels, further reducing the number of experiments.

607 Currently, *doepipeline* leverages cloud computing capability through the Slurm workload

608 managing system. Given the recent development and consolidation of workflow managing

609 systems [44] it would be possible to integrate *doepipeline* with for example SnakeMake [45]

610 or NextFlow [46], similar to other implementations [47,48].

611 During development and testing of *doepipeline* we saw the design space moving back and

612 forth between iterations in the optimization phase. We hypothesized that this behavior was

613 because either the underlying function was not modeled properly or the function was flat

614 within the investigated design space. To counteract this phenomenon we implemented three

615 features; i) no prediction of the optimal factor combination if the predictive power (Q2) of the

616 model was low (default: Q2<0.5), ii) validation of the predicted optimal factor combination

617 by carrying out the pipeline with those factor settings, and iii) shrinking the span of the

618 factors between iterations. After implementing these three features, *doepipeline* consistently

619 converged to satisfactory results.

620 Specifying the pipeline in a YAML file allows for flexible configurations of commands to be

621 run, essentially enabling optimization of any pipeline run on the command line. However, the

622 number of parameters will typically increase with the length of the pipeline under

623 investigation. At the same time there is a soft constraint on the number of parameters that can

624 be investigated simultaneously. This constraint will be related to the problem currently under

625 investigation and depends on the computational complexity of the pipeline, and on the

626 available computational and time resources. Instead of doing a global optimization of

627 parameters, i.e. optimizing the entire pipeline at once, an alternative approach is to run

628 sequential optimizations in which only a section of the pipeline at a time is optimized while

629    keeping the default parameter values for the rest of the pipeline [9]. This type of sequential

630    optimization is not yet fully implemented in *doepipeline* and is a feature for future updates.

631    Sequential optimization of a pipeline currently requires that an optimization is carried out for

632    each step of the pipeline and that the optimized parameter values so obtained are manually

633    updated for the subsequent steps of the pipeline.

## 634    Conclusion

635    Our proposed strategy represents a systematic approach to the optimization of software

636    parameters. Our implementation in the software-agnostic and user-friendly package

637    *doepipeline* could potentially serve as a starting point for experimenters and

638    bioinformaticians who currently rely on default settings or common practice when running

639    their data processing pipelines.

## 640    Declarations

641    **Ethics approval and consent to participate**

642    Not applicable

643    **Availability of data and material**

644    The datasets generated and analyzed during the current study (cases 1-3) are available in

645    NCBI BioProject accessions PRJNA510899 (SRA run: SRR9290761), PRJNA16087 (SRA

646    run: SRR518502), PRJNA548675 (SRA run: SRR9290851). The assembly used as input to

647    case 2 is available at the *doepipeline* github (https://github.com/clicumu/doepipeline). The

648    sequence data used for case 4 is available at the ENA with study accession ERP001960 (run:

649    ERR194147), and the truth callset along with confident regions BED file is available at the

650    Platinum Genomes github repository [29]. The GATK bundle is available at the GATK FTP

651    ftp://ftp.broadinstitute.org/bundle/hg19/.

30

652 **Consent for publication**

653 Not applicable

654 **Competing interests**

655 The authors declare that they have no competing interests.

656 **Funding**

657 This work was funded by the Knut and Alice Wallenberg Foundation (DSv, JT) (2011.0042),

658 the Swedish Research Council (DSv, RS, JT) (2016☐04376), Sartorius AG (RS, JT), the

659 Swedish National Strategic e-Science Research Program eSSENCE (RS, JT), the Swedish

660 Ministry of Defence (DSu, AS) (FOI project A404018) and the Swedish Civil Contingencies

661 Agency (DSu, AS) (FOI project B4662 B2Forensics).

662 **Authors' contributions**

663 DSv, AS and JT initiated the study, DSv and RS designed the algorithm and implemented the

664 python code. DSv and DSu analyzed and interpreted the four example cases. All authors

665 wrote, read and approved the final manuscript.

## 666 Acknowledgements

## 675 References

676   1.   DePristo MA, Banks E, Poplin R, Garimella K V, Maguire JR, Hartl C, et al. A framework for

677       variation discovery and genotyping using next-generation DNA sequencing data. Nat Genet

678       [Internet].   2011   May   [cited   2018   Jan   17];43(5):491–8.   Available   from:

679       http://www.ncbi.nlm.nih.gov/pubmed/21478889

680   2.   Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, Del Angel G, Levy-Moonshine A, et al.

681       From FastQ data to high confidence variant calls: the Genome Analysis Toolkit best practices

682       pipeline. Curr Protoc Bioinforma [Internet]. 2013 [cited 2018 Jan 17];43(1110):11.10.1-33.

683       Available from: http://www.ncbi.nlm.nih.gov/pubmed/25431634

684   3.   Li H. Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics [Internet]. 2018

685       [cited   2018   Dec   20];34(18):3094–100.   Available   from:

686       http://www.ncbi.nlm.nih.gov/pubmed/29750242

687   4.   Fisher RA. The design of experiments. Edinburgh/London: Oliver and Boyd; 1935.

688   5.   Eriksson L, Johansson E, Kettaneh-Wold N, Wikström C, Wold S. Design of experiments :

689       principles and applications [Internet]. Umeå: Umetrics Academy; 2008. Available from:

690       http://www.umetrics.com

691   6.   Box GEP, Hunter WG, Hunter JS. Statistics for experimenters□: an introduction to design,

692        data analysis, and model building. New York: Wiley; 1978. (Wiley series in probability and

693        mathematical statistics, 0277-2728).

694   7.   Dismuke C, Lindrooth R. Ordinary least squares. Methods Des Outcomes Res. 2006;93:93–

695        104.

696   8.   Surowiec I, Vikström L, Hector G, Johansson E, Vikström C, Trygg J. Generalized Subset

697        Designs in Analytical Chemistry. Anal Chem. 2017;89(12):6491–7.

698   9.   Eliasson M, Rännar S, Madsen R, Donten MA, Marsden-Edwards E, Moritz T, et al. Strategy

699        for Optimizing LC-MS Data Processing in Metabolomics: A Design of Experiments

700        Approach. Anal Chem [Internet]. 2012 Aug 7 [cited 2019 Apr 18];84(15):6869–76. Available

701        from: http://www.ncbi.nlm.nih.gov/pubmed/22823568

702   10.  Derringer G, Suich R. Simultaneous Optimization of Several Response Variables. J Qual

703        Technol [Internet]. 1980 Oct 22 [cited 2018 Mar 2];12(4):214–9. Available from:

704        https://www.tandfonline.com/doi/full/10.1080/00224065.1980.11980968

705   11.  Svensson K, Sjödin A, Byström M, Granberg M, Brittnacher MJ, Rohmer L, et al. Genome

706        sequence of Francisella tularensis subspecies holarctica strain FSC200, isolated from a child

707        with tularemia. J Bacteriol [Internet]. 2012 Dec [cited 2018 Dec 19];194(24):6965–6.

708        Available from: http://www.ncbi.nlm.nih.gov/pubmed/23209222

709   12.  seqkt [Internet]. Available from: https://github.com/lh3/seqtk

710   13.  Grüning B, Dale R, Sjödin A, Chapman BA, Rowe J, Tomkins-Tinch CH, et al. Bioconda:

711        sustainable and comprehensive software distribution for the life sciences. Nat Methods

712        [Internet]. 2018 Jul 2 [cited 2018 Dec 20];15(7):475–6. Available from:

713        http://www.nature.com/articles/s41592-018-0046-7

714   14.  Simpson JT, Wong K, Jackman SD, Schein JE, Jones SJ, Birol I. ABySS: A parallel assembler

715        for short read sequence data. [cited 2018 Jun 14]; Available from: www.genome.org.

716   15.  Jackman SD, Vandervalk BP, Mohamadi H, Chu J, Yeo S, Hammond SA, et al. ABySS 2.0:

717        resource-efficient assembly of large genomes using a Bloom filter. Genome Res [Internet].

718        2017        [cited        2018        Dec        19];27(5):768–77.        Available        from:

33

719        http://www.ncbi.nlm.nih.gov/pubmed/28232478

720   16.   Earl D, Bradnam K, St John J, Darling A, Lin D, Fass J, et al. Assemblathon 1: a competitive

721        assessment of de novo short read assembly methods. Genome Res [Internet]. 2011 Dec 1 [cited

722        2018 Dec 20];21(12):2224–41. Available from:

723        http://www.ncbi.nlm.nih.gov/pubmed/21926179

724   17.   Bradnam KR, Fass JN, Alexandrov A, Baranay P, Bechner M, Birol İİ, et al. Assemblathon 2:

725        evaluating de novo methods of genome assembly in three vertebrate species. Gigascience

726        [Internet]. 2013 [cited 2018 Dec 12];2(1):10. Available from: http://arxiv.org/abs/1301.5406

727   18.   Fastaq [Internet]. Available from: https://github.com/sanger-pathogens/Fastaq

728   19.   seqstats [Internet]. Available from: https://github.com/clwgg/seqstats

729   20.   Boetzer M, Pirovano W. SSPACE-LongRead: scaffolding bacterial draft genomes using long

730        read sequence information. BMC Bioinformatics [Internet]. 2014 Jun 20 [cited 2018 Jul

731        27];15(1):211. Available from:

732        http://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-15-211

733   21.   Breitwieser FP, Lu J, Salzberg SL. A review of methods and databases for metagenomic

734        classification and assembly. Brief Bioinform [Internet]. 2017 Sep 23 [cited 2018 Dec 20];

735        Available from: http://www.ncbi.nlm.nih.gov/pubmed/29028872

736   22.   Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact

737        alignments. Genome Biol [Internet]. 2014 Mar 3 [cited 2018 Dec 19];15(3):R46. Available

738        from: http://genomebiology.biomedcentral.com/articles/10.1186/gb-2014-15-3-r46

739   23.   Breitwieser FP, Baker DN, Salzberg SL. KrakenUniq: confident and fast metagenomics

740        classification using unique k-mer counts. Genome Biol [Internet]. 2018 Dec 16 [cited 2018

741        Dec 20];19(1):198. Available from:

742        https://genomebiology.biomedcentral.com/articles/10.1186/s13059-018-1568-0

743   24.   Supernat A, Vidarsson OV, Steen VM, Stokowy T. Comparison of three variant callers for

744        human whole genome sequencing. Sci Rep [Internet]. 2018 Dec 14 [cited 2019 May

745        9];8(1):17851. Available from: http://www.ncbi.nlm.nih.gov/pubmed/30552369

746   25.   Eberle MA, Fritzilas E, Krusche P, Källberg M, Moore BL, Bekritsky MA, et al. A reference

34

747     data set of 5.4 million phased human variants validated by genetic inheritance from sequencing

748     a three-generation 17-member pedigree. Genome Res. 2017;27(1):157–64.

749  26. Zook JM, Chapman B, Wang J, Mittelman D, Hofmann O, Hide W, et al. Integrating human

750     sequence data sets provides a resource of benchmark SNP and indel genotype calls. Nat

751     Biotechnol [Internet]. 2014 Mar [cited 2014 Jul 19];32(3):246–51. Available from:

752     http://www.ncbi.nlm.nih.gov/pubmed/24531798

753  27. Zook JM, McDaniel J, Parikh H, Heaton H, Irvine SA, Trigg L, et al. Reproducible integration

754     of multiple sequencing datasets to form high-confidence SNP, indel, and reference calls for

755     five human genome reference materials. bioRxiv [Internet]. 2018 Mar 13 [cited 2019 May

756     8];281006. Available from: https://www.biorxiv.org/content/10.1101/281006v1

757  28. Krusche P, Trigg L, Boutros PC, Mason CE, Vega FMD La, Moore BL, et al. Best Practices

758     for Benchmarking Germline Small Variant Calls in Human Genomes. bioRxiv [Internet]. 2018

759     Feb     23     [cited     2019     May     8];270157.     Available     from:

760     https://www.biorxiv.org/content/10.1101/270157v1.full

761  29. Platinum Genomes GitHub repository / hg19 hybrid truth set [Internet]. Available from:

762     https://illumina.github.io/PlatinumGenomes/?prefix=2017-1.0/hg19/hybrid

763  30. Picard [Internet]. Available from: http://broadinstitute.github.io/picard

764  31. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.

765     Bioinformatics [Internet]. 2009 Jul 15 [cited 2018 Jul 5];25(14):1754–60. Available from:

766     http://www.ncbi.nlm.nih.gov/pubmed/19451168

767  32. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. 2013

768     Mar 16 [cited 2019 May 8]; Available from: http://arxiv.org/abs/1303.3997

769  33. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The Genome

770     Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing

771     data. Genome Res [Internet]. 2010 Sep [cited 2018 Jul 5];20(9):1297–303. Available from:

772     http://www.ncbi.nlm.nih.gov/pubmed/20644199

773  34. Krusche    P.    Haplotype    comparison    tools    /    hap.py    [Internet].    Available    from:

774     http://github.com/illumina/hap.py

35

775    35.    conda-forge [Internet]. Available from: https://conda-forge.org/

776    36.    doepipeline      (conda-forge)      [Internet].      Available      from:      https://anaconda.org/conda-

777           forge/doepipeline

778    37.    PyDOE2 [Internet]. Available from: https://github.com/clicumu/pyDOE2

779    38.    Yoo AB, Jette MA, Grondona M. SLURM: Simple Linux Utility for Resource Management.

780           In  Springer,  Berlin,  Heidelberg;  2003  [cited  2018  Dec  19].  p.  44–60.  Available  from:

781           http://link.springer.com/10.1007/10968987_3

782    39.    VelvetOptimizer [Internet]. Available from: https://github.com/tseemann/VelvetOptimiser

783    40.    Zerbino DR, Birney E. Velvet: algorithms for de novo short read assembly using de Bruijn

784           graphs. Genome Res [Internet]. 2008 May 1 [cited 2018 Dec 20];18(5):821–9. Available from:

785           http://www.ncbi.nlm.nih.gov/pubmed/18349386

786    41.    Chikhi R, Medvedev P. Informed and automated k-mer size selection for genome assembly.

787           Bioinformatics  [Internet].  2014  Jan  1  [cited  2018  Dec  12];30(1):31–7.  Available  from:

788           https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btt310

789    42.    Bergstra J, Bengio Y. Random search for hyper-parameter optimization. J Mach Learn Res

790           [Internet].  2012;13:281–305.  Available  from:  papers3://publication/uuid/1190E1AB-0319-

791           40C5-81CD-7207784965DE

792    43.    Snoek  J,  Larochelle  H,  Adams  RP.  Practical  Bayesian  Optimization  of  Machine  Learning

793           Algorithms. Adv Neural Inf Process Syst [Internet]. 2012 Jun 13 [cited 2019 Jun 6]; Available

794           from: http://arxiv.org/abs/1206.2944

795    44.    Karim MR, Michel A, Zappa A, Baranov P, Sahay R, Rebholz-Schuhmann D. Improving data

796           workflow systems with cloud services and use of open data for bioinformatics research. Brief

797           Bioinform  [Internet].  2018  Sep  28  [cited  2019  Jun  20];19(5):1035–50.  Available  from:

798           https://academic.oup.com/bib/article/19/5/1035/3737318

799    45.    Koster J, Rahmann S. Snakemake--a scalable bioinformatics workflow engine. Bioinformatics

800           [Internet].   2012   Oct   1   [cited   2019   Jun   20];28(19):2520–2.   Available   from:

801           https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/bts480

802    46.    Di Tommaso P, Chatzou M, Floden EW, Barja PP, Palumbo E, Notredame C. Nextflow

36

803    enables reproducible computational workflows. Nat Biotechnol [Internet]. 2017 Apr 1 [cited

804    2019 Jun 20];35(4):316–9. Available from: http://www.nature.com/articles/nbt.3820

805    47.    Holl S, Mohammed Y, Zimmermann O, Palmblad M. Scientific workflow optimization for

806    improved peptide and protein identification. BMC Bioinformatics [Internet]. 2015 Dec 3 [cited

807    2019          Jun          20];16(1):284.          Available          from:

808    http://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-015-0714-x

809    48.    Palmblad M, Lamprecht A-L, Ison J, Schwämmle V. Automated workflow composition in

810    mass spectrometry-based proteomics. Wren J, editor. Bioinformatics [Internet]. 2019 Feb 15

811    [cited          2019          Jun          20];35(4):656–64.          Available          from:

812    https://academic.oup.com/bioinformatics/article/35/4/656/5060940

813

## 814    Additional material

815    **Additional file 1**

816    Excel file (.xlsx).

817    Complete experimental sheets for case 1, *doepipeline* iterations 1 and 2. Contains factor

818    settings and response values for all experiments in these iterations.

819    **Additional file 2**

820    Excel file (.xlsx).

821    Complete experimental sheet for case 1, grid search. Contains factor settings and response

822    values for all executed experiments.

823    **Additional file 3**

824    Excel file (.xlsx).

825    Complete experimental sheets for case 2, *doepipeline* iterations 1, 2, and 3. Contains factor

826    settings and response values for all experiments in these iterations.

827    **Additional file 4**

828    Excel file (.xlsx).

829    Complete experimental sheets for case 2, grid search. Contains factor settings and response

830    values for all executed experiments.

831    **Additional file 5**

832    Excel file (.xlsx).

833    Complete experimental sheets for case 3, *doepipeline* iterations 1, 2, and 3. Contains factor

834    settings and response values for all experiments in these iterations.

835    **Additional file 6**

836    Excel file (.xlsx).

837    Complete experimental sheet for case 3, grid search. Contains factor settings and response

838    values for all executed experiments.

839    **Additional file 7**

840    Excel file (.xlsx).

841    Complete experimental sheets for case 4, step 1, *doepipeline* iterations 1, 2, and 3. Contains

842    factor settings and response values for all experiments in these iterations.

843    **Additional file 8**

844    Excel file (.xlsx).

845    Complete experimental sheets for case 4, step 2, *doepipeline* iterations 1, 2, 3, and 4.

846    Contains factor settings and response values for all experiments in these iterations.

847    **Additional file 9**

848    Excel file (.xlsx).

849    Complete experimental sheet for case 4, step 1, grid search. Contains factor settings and

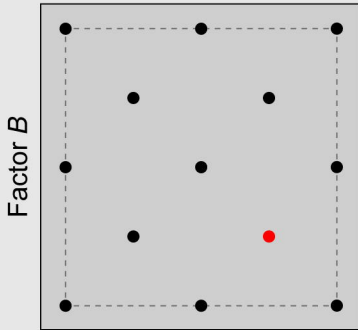850    response values for all executed experiments.

851    **Additional file 10**

852    Excel file (.xlsx).

853     Complete experimental sheet for case 4, step 2, grid search. Contains factor settings and

854     response values for all executed experiments.
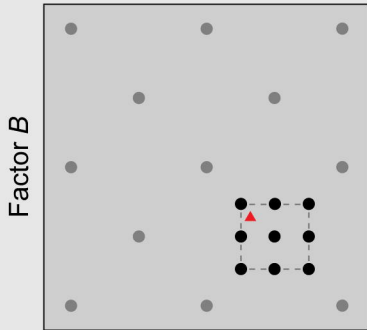
855

**(a)**

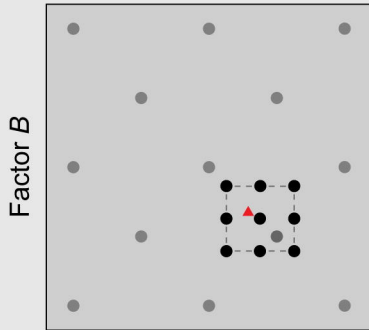Screening phase

iteration 1

Factor $B$

Factor $A$

**(b)**

Optimization phase

iteration 2

Factor $B$

Factor $A$

iteration 3

Factor $B$

Factor $A$