# TranSPHIRE: Automated and feedback-optimized on-the-fly processing for cryo-EM

Markus Stabrin, Fabian Schoenfeld, Thorsten Wagner, Sabrina Pospich,

Christos Gatsogiannis, Stefan Raunser*

*Department of Structural Biochemistry, Max Planck Institute of Molecular Physiology,*

*Otto-Hahn-Straße 11, 44227 Dortmund, Germany.*

(*corresponding author: stefan.raunser@mpi-dortmund.mpg.de)

## Abstract

Single particle electron cryomicroscopy (cryo-EM) requires full automation to allow high-throughput structure determination which is especially important for drug discovery research. Although several software packages exist where parts of the cryo-EM pipeline are automated, a complete solution that offers reliable, quality-optimized on-the-fly processing, resulting in a high-resolution three-dimensional reconstruction does not exist. Here we present TranSPHIRE: A software package for fully automated processing of cryo-EM data sets during data acquisition. TranSPHIRE transfers data from the microscope, automatically applies the common pre-processing steps, picks particles, performs 2D clustering, and 3D refinement parallel to image recording. Importantly, TranSPHIRE introduces a machine learning-based feedback loop to re-train its internally used picking model to adapt to any given data set live during processing. This elegant approach enables TranSPHIRE to process data more effectively, producing high-quality particle stacks. TranSPHIRE collects, and displays all microscope settings and metrics generated by its individual tools, in order to allow users to quickly evaluate data during acquisition. TranSPHIRE can run on a single work station and also includes the automated processing of filaments.

## Introduction

Single particle electron cryomicroscopy (cryo-EM) has successfully established itself as a prime method to determine the three-dimensional structure of macromolecular complexes at close to atomic resolution [1,2]. The technique has therefore the potential to become a key tool for drug discovery research [3]. However, single particle analysis (SPA) studies still require large amounts of processing time, expert knowledge, and computational resources. With the number of modern high-throughput microscopes growing rapidly, there is an urgent demand for a robust, automated processing pipeline that requires little to no user intervention. This need is felt especially in the field of drug discovery [3].

In many cases, data sets that were recorded for several days and can include 10,000 to 20,000 movies turn out to be unusable for high-resolution structure determination during subsequent data processing. It is therefore necessary for users to obtain feedback on the quality of their data immediately during recording. This enables them to decide whether or not to continue a session, adjust any of the acquisition parameters at the microscope, and compare different grids. This can only be achieved when processing the data in parallel to data acquisition. A fully automated pipeline requires streamlined data transfer and automated pre-processing and processing workflows, free of any user bias.

Although several software packages partially address these issues [4-10], a complete solution that offers reliable, quality-optimized and flexible on-the-fly processing during data acquisition resulting in a high-resolution 3D reconstruction does not yet exist. CryoFLARE [4], for example, performs live analysis and processing parallel to data acquisition, but only to the level of 2D classification and lacks the ability to perform *ab initio* 3D reconstructions or high resolution refinements. Similarly, Focus [8], Appion [10], and Warp [6] do not produce 2D class averages and 3D

53    reconstructions. The latter two are less flexible than other offerings by being

54    restricted to either collecting data with Leginon [11] in case of Appion or exclusive

55    compatibility with Windows and Warp-native tools. All three software packages

56    concentrate on data acquisition and associated parameters but not on the

57    optimization of data processing which is an important prerequisite for automated

58    structure determination. The non-interactive data pre-processing in Relion-3 [5] offers,

59    similar to Focus [8], some flexibility in terms of tool integration, but hinders the

60    implementation of more complicated cryo-EM processing by making advanced

61    parameters only accessible via manual scripting, rather than its GUI. Both Relionit [5]

62    and Scipion [7] share the same accessibility issue of quality metrics, where no values

63    are automatically plotted and updated during processing. Instead, the user has to

64    step in and trigger the compilation of a log-file that contains a mix of metrics for all

65    processed data; any specific values of an individual micrograph have to be found

66    manually. This makes assessment problematic, especially for beginners in the field.

67    Here we present TranSPHIRE, a fully automated pipeline for on-the-fly

68    processing of cryo-EM data. It combines deep learning tools with a novel, feedback-

69    driven approach to re-train the integrated crYOLO particle picker [12] during ongoing

70    pre-processing. This allows TranSPHIRE to perform GPU accelerated 2D classification

71    to provide high-quality 2D class averages and, subsequently, 3D reconstructions from

72    clean data. This gives experimentalists the means to quickly evaluate both the quality

73    of their data sets as well as their chosen microscope settings during data acquisition.

74    A combination of new and improved tools allows TranSPHIRE to provide users with

75    the strongest early results in the shortest amount of time, without the need to

76    outsource the computational load to a computer cluster, or the need for user

77    intervention. Importantly, it allows users to perform automated high-throughput on-

78    the-fly screenings for different buffer conditions or ligands of interest as well as to

3

79   fine-tune the workflow for the respective target-protein and perform digital
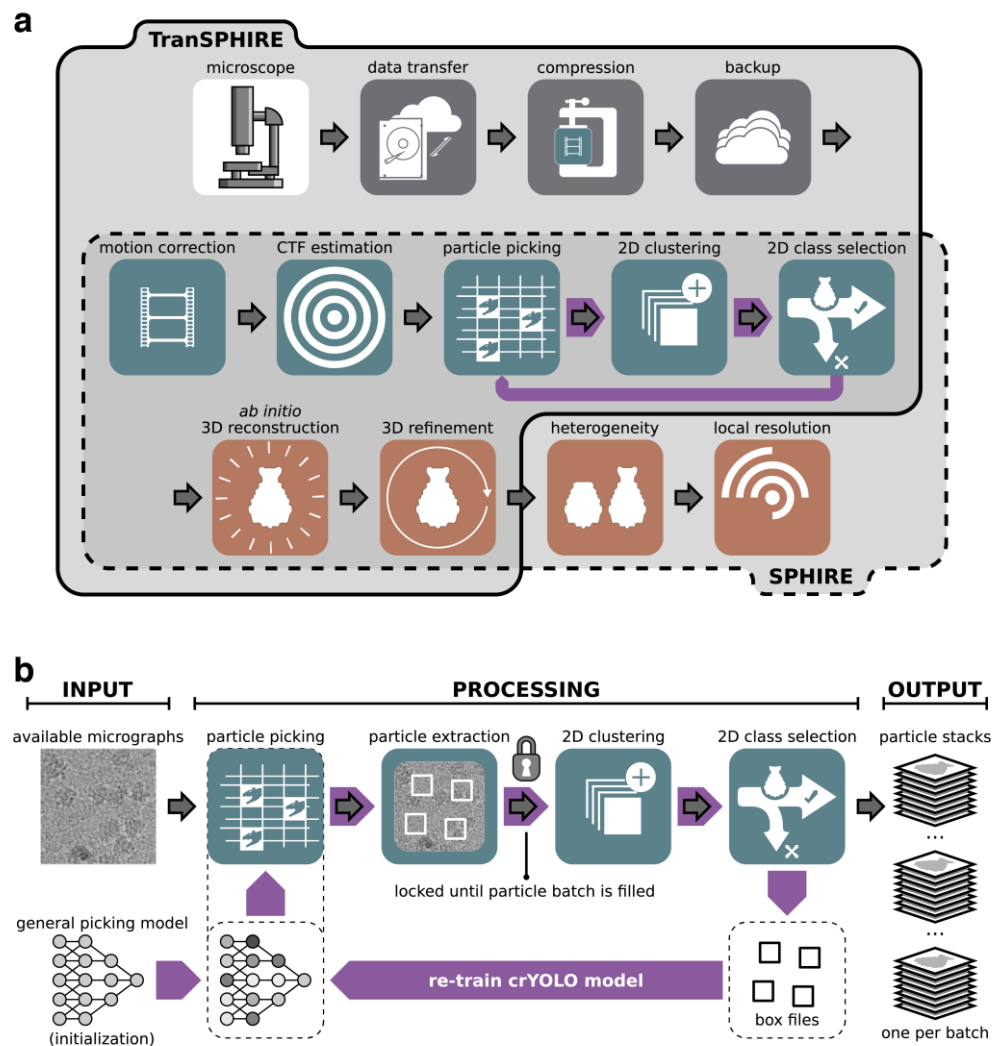
80   purification during image acquisition.

81

## Results

**General setup, functionality and layout of TranSPHIRE**

84   TranSPHIRE is an automated pipeline for processing cryo-EM data sets (Figure 1). It is

85   developed in Python 3 to run on Linux, and is available online for free. TranSPHIRE

86   performs parallelized data transfer and flexibly integrates a range of commonly used

87   pre-processing tools, as well as the advanced processing tools of the SPHIRE package

88   [13]. Using these tools, TranSPHIRE implements a fully automated pipeline to process

89   cryo-EM data on-the-fly during data acquisition. TranSPHIRE is designed to allow

90   users to make the best use of their available resources by prioritizing data analysis,

91   presenting early results, and using machine learning tools to identify and process

92   only those parts of the data that contribute to high quality results.

93   TranSPHIRE is controlled via an easy-to-use GUI that allows users to set up a

94   session, and choose and configure the desired tools to use (Supplementary Figure 1).

95   For pre-processing, the TranSPHIRE pipeline integrates MotionCor2 [14] and Unblur [15]

96   for beam induced motion correction with dose weighting; as well as CTFFIND4 [16],

97   CTER [17], and GCTF [18] for CTF estimation. This modularized integration is entirely

98   parameterized, allowing experimentalists to both choose their preferred tools as well

99   as configure them as needed – all without leaving the TranSPHIRE GUI. Available

100  parameters are sorted by level of usage ("main", "advanced", and "rare") to highlight

101  and help identify the most commonly adjusted parameters for each tool.

102  During the session, TranSPHIRE automatically parallelizes the batch-wise

103  processing of incoming micrographs, outsources computationally expensive steps to

4

104

**Figure 1. The TranSPHIRE pipeline and the SPHIRE backend. (a) Upper register (solid line):** *Overview of the integrated TranSPHIRE pipeline and all automated processing steps. The pipeline includes file management tasks, i.e., parallelized data transfer, file compression, and file backup (grey); 2D processing, i.e., motion correction, CTF estimation, particle picking, 2D clustering, and 2D class selection (turquoise); and 3D processing, i.e., ab initio 3D reconstruction and 3D refinement (red). Additionally, the pipeline includes an automated feedback loop optimization to adapt picking to the current data set during runtime (purple).* **Lower register (dotted line):** *The SPHIRE software package forms the backend for TranSPHIRE and offers the tools used for 2D and 3D processing. SPHIRE includes additional tools for advanced processing, such as heterogeneity analysis and local resolution determination.* **(b)** *The TranSPHIRE feedback loop. Grey arrows indicate the flow of data processing. Red arrows indicate the flow of the feedback loop.* **Left (input):** *Micrographs are initially picked using the crYOLO general model.* **Center (processing):** *Particles are picked and extracted. Once*

119    *a pre-defined number of particles have been accumulated, the pipeline performs 2D*

120    *classification; the resulting 2D class averages are labeled as either "good" or "bad" by*

121    *Cinderella. Class labels and crYOLO box files are then used to re-train crYOLO and*

122    *adapt its internal model to the processed data. In the next feedback round this*

123    *updated model is used to re-pick the data.* **Right (output):** *After five feedback rounds,*

124    *the complete data set is picked with the final optimized picking model and 2D*

125    *classified in batches. For every batch a particles stack of "good" particles is created*

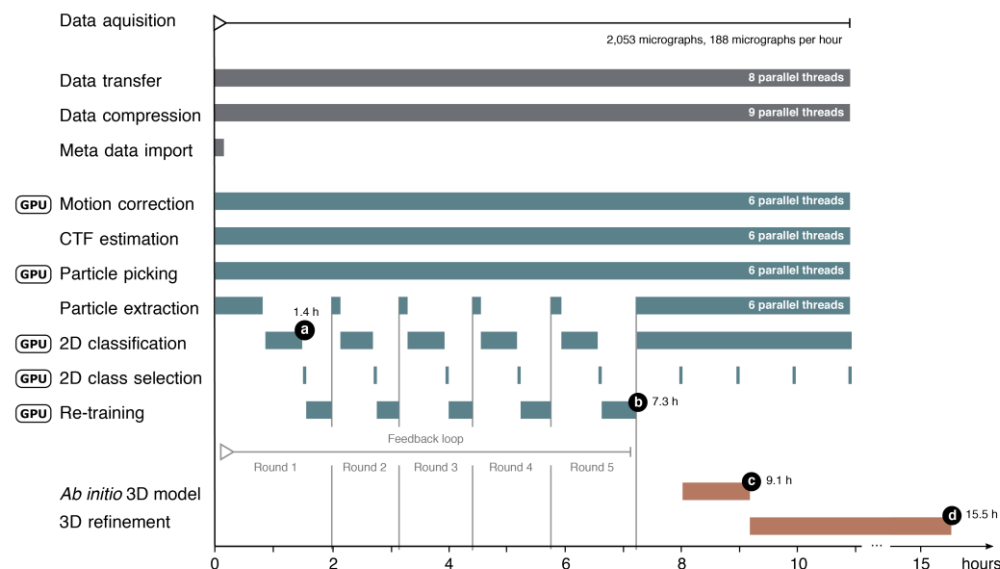126    *and available for 3D processing.*

127

128    available GPUs, and produces preliminary 2D class averages and 3D reconstructions

129    based on the most recently processed batch of data (Figure 2, Supplementary Figure

130    2). Through the optimal distribution of processes, TranSPHIRE runs on-the-fly for a

131    wide range of data acquisition settings using a single workstation (Supplementary

132    Figure 3; see Methods for details about hardware). Moreover, TranSPHIRE can catch-

133    up with the speed of the acquisition after the initial delay due to the feedback loop

134    (see below) for routinely used data acquisition schemes (Figure 3c). Thus, initial 2D

135    class averages and 3D reconstructions are available within a few hours after starting

136    the data collection (Figure 2, Supplementary Figure 3).

137        Throughout the processing, TranSPHIRE collects all data quality metrics

138    produced by its individual tools, links them with the relevant micrographs where

139    appropriate, and presents them front and center in its GUI (Supplementary Figure 1).

140        Optionally, notifications for early milestones such as 2D class averages and

141    preliminary 3D maps, can also be sent via email. These features enable

142    experimentalists to both identify and address any issues as soon as they surface

143    during data acquisition, without requiring constant user supervision. Additionally, all

144    results produced by the integrated tools during processing are also copied in parallel

145    to the pre-defined workstation and backup locations (Supplementary Figure 2). To

146    support interoperability with existing packages, all pre-processing steps until particle

6

147     picking support the file formats used in both SPHIRE and RELION; for later processing

148     steps utilities are included to convert SPHIRE files into RELION .star.

149



150

151     ***Figure 2. Timeline of the TranSPHIRE pipeline.*** *Timeline depicting the parallel*

152     *execution of the processes of the TranSPHIRE pipeline. Timings are based on a Tc*

153     *holotoxin data set consisting of 2,053 micrographs, each containing 36 particles on*

154     *average, collected at a speed of 188 micrographs per hour (K2 super-resolution, 40*

155     *frames). TranSPHIRE ran on-the-fly up to the creation of an ab initio 3D*

156     *reconstruction using default settings. Important milestones are denoted in black:* ***(a)***

157     *first 2D class averages produced after 1.4h;* ***(b)*** *end of the feedback loop after 7.3h;*

158     ***(c)*** *ab initio 3D reconstruction after 9.1h; and* ***(d)*** *final 3D reconstruction of the first*

159     *batch of particles after 15.5h.*

160

161     **Transfer and pre-processing**

162     Once a session starts, TranSPHIRE automatically detects and transfers new

163     micrographs from the camera computer of the microscope (Figure 1, Supplementary

164     Figure 2). These data are moved in parallel to several, user-specified locations e.g. a

165     work station or cluster for processing, and a backup storage server. In case of the

166     latter, TranSPHIRE also automatically compresses the data to preserve storage space.

167     Copy locations may also include additional spaces such as transportable hard discs. If

7

168    desired, TranSPHIRE further renames files, and deletes images from the camera

169    workstation in order to free up more space to enable continuous data collection. It

170    also extracts meta data such as acquisition time, grid square, hole number and

171    coordinates, spot scan, and phase plate position from .xml files provided by EPU or

172    .gtg files provided by Latitude S.

173        During the ongoing data transfer, any data that has already been copied is

174    pre-processed in parallel (Figure 2, Supplementary Figure 2). During setup, users can

175    choose to perform motion correction using either MotionCor2 [14] or Unblur[15]. While

176    motion correction is performed, TranSPHIRE presents all relevant metrics, such as the

177    average shift per frame, or the overall shift per micrograph (Supplementary Figure 1).

178    For CTF estimation, users can set up TranSPHIRE to use either CTFFIND4 [16], CTER [17],

179    and GCTF [18]. Depending on whether or not CTF estimation on movies is activated in

180    TranSPHIRE, CTF estimation is performed in parallel to motion correction

181    (Supplementary Figure 2). The metrics extracted and displayed by TranSPHIRE

182    include defocus, astigmatism, and the resolution limit (Supplementary Figure 1).

183    Combined with the information gathered during motion correction, these values

184    allow experimentalists to assess the performance and alignment of the microscope

185    during acquisition, and adjust any thresholds to automatically discard low quality

186    micrographs as necessary.

187        For particle picking the TranSPHIRE pipeline integrates crYOLO [12], our state of

188    the art deep learning particle picker. During picking, TranSPHIRE displays the particles

189    picked per micrograph, which allows users to assess the picking performance and

190    overall sample quality (Supplementary Figure 1).

191        Once a fixed threshold of picked particles is reached (Supplementary Figure

192    4; also see Methods), TranSPHIRE launches 2D classification using a GPU accelerated

193    version of ISAC2 [19] (Figure 1). ISAC2 limits the number of class members to spread the

194    given particles across multiple classes which prevents individual classes from growing

195    too large. This results in sharp, equal-sized, and reproducible classes that contain all

196    possible orientations exceeding the minimum class size. They enable

197    experimentalists to reliably assess particle orientations and overall quality, and help

198    to identify possible issues such as preferred orientations or heterogeneity.

199    The 2D class averages are then sorted by Cinderella [20], our integrated deep

200    learning tool for 2D class selection. Cinderella labels the given 2D classes as either

201    "good" or "bad" and determines which class averages and, thereby, particles are

202    used for further processing. This results in an automatic cleaning of the data and

203    allows TranSPHIRE to process only the relevant subset of a given data set, thereby

204    dramatically lowering the amount of data processed by the computationally

205    expensive steps of 3D reconstruction and refinement.

206

**207    Optimizing particle picking using a machine learning-fueled feedback loop**

208    For any cryo-EM pipeline the ability to reliably perform high quality picking

209    irrespective of the data at hand is essential. This poses a challenge when processing

210    is to be automated, as this immediately excludes any user intervention such as

211    manual inspection of the picking results. The latter is especially relevant if a sample is

212    unknown to the picking procedure, or is otherwise difficult to process, e.g. due to

213    contamination or interfering conformational states – issues that usually need to be

214    identified by a qualified expert before processing can continue.

215    TranSPHIRE solves these issues by introducing a machine learning based

216    feedback loop that repeatedly re-trains the fully integrated crYOLO [12] deep learning

217    particle picker during data acquisition to adapt picking to the given data set (Figure

218    1b). This enables crYOLO to specifically target those particles that end up in stable 2D

219    class averages, while, at the same time, learning to disregard particles that do not.

220 First, incoming motion corrected micrographs are forwarded to crYOLO for picking.

221 Once a batch of 20,000 picked particles has been accumulated, it is handed over to

222 our GPU accelerated version of the 2D classification algorithm ISAC2 [19]

223 (Supplementary Figure 3). Here we determine which particles can be used to create

224 stable 2D class averages, and reject the particles that cannot be accounted for. The

225 newly produced 2D class averages are given to our deep learning tool Cinderella [20],

226 which labels each class average as either "good" or "bad". At this point, the particles

227 of the "good" classes are used to re-train crYOLO and update its internal model.

228 Specifically, we randomly select a maximum of 50 micrographs that contain particles

229 that ended up in the "good" classes for the re-training (for details see Methods).

230 Once the training and thus the first feedback round has completed, processing re-

231 starts using the optimized picking model (Figure 1b).

232 The TranSPHIRE feedback loop iterates five times, which has proven sufficient

233 to achieve convergence in our experiments, and afterwards is not repeated for the

234 remainder of the data acquisition. Re-training crYOLO [12] to become increasingly more

235 proficient at targeting particles that end up in "good" classes has the additional

236 benefit of trimming down the overall size of the data set. Though the pre-processing

237 of cryo-EM data is already time consuming, the following 3D refinement requires

238 even more computational power. While it is usually customary to process as much

239 data as possible, the computational cost of 3D refinement usually does not scale

240 linearly, and such an approach will not be sustainable in the near future. This is

241 further exacerbated by the fact that image acquisition speeds and sizes of data sets

242 are both growing rapidly. Because of this, the aim should be to process as little data

243 as necessary, without harming the quality of the final reconstruction. Fortunately, it

244 is known that cryo-EM data sets contain a large amount of unusable data that can be

245 safely discarded – if we have a way to reliably ensure that we keep those data that

246    we are actually interested in. The TranSPHIRE feedback loop offers this functionality

247    and provides quality in quantity.

248

249    ***Ab initio* 3D model reconstruction and 3D refinement**

250    To compute a 3D reconstruction, the particles included in all classes labeled "good"

251    by Cinderella are extracted and form a clean, high-resolution particle stack. If there is

252    no initial 3D reference provided to TranSPHIRE, the pipeline waits until at least 200

253    (by default) "good" classes have been accumulated. The respective 2D class averages

254    are then used to create a reproducible, *ab initio* 3D reconstruction using SPHIRE

255    RVIPER [13,21] (Figure 1). This provides a first view of the structure of the target protein

256    and a first impression of the conformational state.

257         The initial 3D reference is then used by TranSPHIRE to initialize the 3D

258    refinement using SPHIRE MERIDIEN (Figure 2). While the initial map is computed only

259    once, a new 3D refinement is started every time another set of 40,000 (by default)

260    "good" particles has been accumulated.

261         Note that in contrast to SPHIRE RVIPER, which only uses the first 2D class

262    averages, SPHIRE MERIDIEN uses all particles subsumed by the last batch of "good"

263    particles. The fully automated creation of an initial 3D map and continuous

264    production of a series of refined reconstructions based on that latest data enables

265    TranSPHIRE to present high-resolution structures already during data acquisition.

266         This enables for a more detailed, on-the-fly evaluation by the user, such as

267    analyzing the conformational state and/or confirming whether and where a ligand is

268    bound. By providing a series of reconstructions – one for every batch, TranSPHIRE

269    also offers a time-resolution of the data set, enabling experimentalists to gauge the

270    quality of their data over time throughout data acquisition.

271    With the following three experiments we illustrate the capabilities of

272  TranSPHIRE to automatically adapt to unknown data, make use of prior knowledge to

273  selectively target the conformational subpopulation within a sample and process
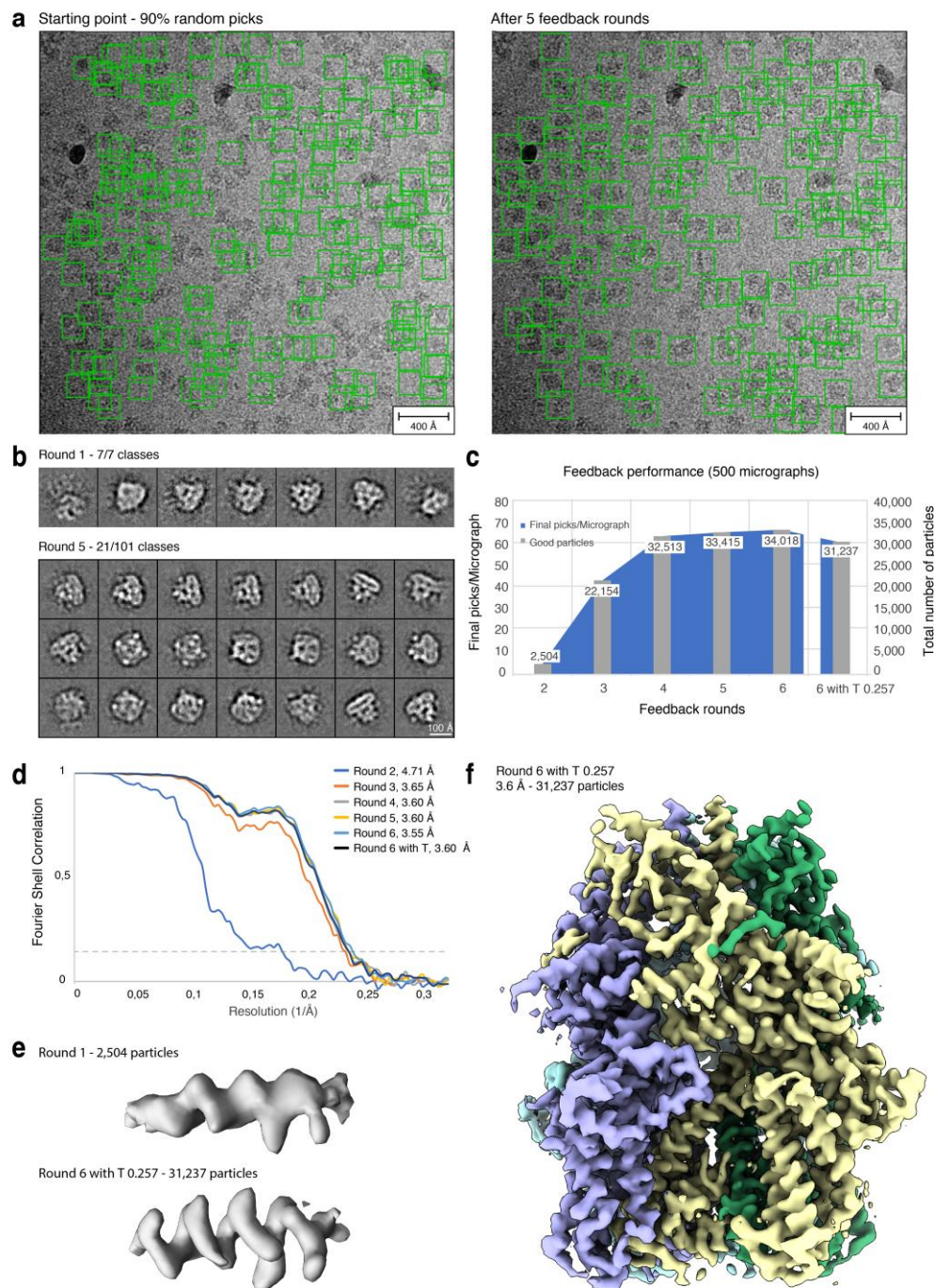
274  filamentous data.

275

276  **Learning to pick an initially unknown membrane channel without user intervention**

277  Similar to crYOLO, many modern particle picking programs are based on machine

278  learning, where an internal model is trained to recognize particles within

279  micrographs [6,22,23]. While this method features an inherent capacity to generalize to

280  unseen data sets, this ability is limited. Therefore, reliable picking can usually not be

281  guaranteed out of the box when samples differ too much from the original training

282  data of the network. Samples might also be of unusually low contrast, or an unknown

283  form of contamination is encountered. While such issues can be overcome by adding

284  the problematic data to the training set, this requires manual user intervention on

285  multiple levels. First, the insufficient picking capability has to be detected; second, an

286  experienced experimentalist has to pick a small amount of training data by hand; and

287  third, the network has to be re-trained manually.

288    The TranSPHIRE feedback loop resolves this issue and entirely foregoes the

289  need for user intervention even when facing data that is either unknown to the

290  picking model or yields insufficient picking results for any other reason. To

291  demonstrate this ability, we processed a data set of the TRPC4 membrane protein

292  channel with the TranSPHIRE feedback loop using a picking model without any prior

293  knowledge of this protein (Figure 3). Specifically, to ensure the sample was unknown

294  to crYOLO at the start of the feedback loop, we removed all four TRP channel data

295  sets normally included in the training data of the crYOLO general model. Additionally,

296  in order to simulate a bad generalization of crYOLO we randomized 90% of all picks in

297



298
299   ***Figure 3. Processing the TRPC4 membrane channel using a deliberately hampered***

300   ***picking model. (a)** To simulate low quality picking, only 10% of the initial crYOLO*

301   *picks were used while the remaining 90% were re-positioned randomly **(left)**. After*

302   *the feedback loop crYOLO reliably picks the TRPC4 particles **(right)**. **(b)** Total amount*

303   *of 2D class averages produced in the first iteration of the feedback loop **(top)** and 21*

304   *representative averages produced in the final iteration of the feedback loop **(bottom)**.*

305   ***(c)** Progression of the number of particles labeled "good" when applying the*

306   *intermediate picking models of the feedback loop to a fixed subset of 500*

307     *microographs. The curve flattens out in the final iterations, indicating the convergence*

308     *of the feedback loop optimization. **(d)** Fourier shell correlation (FSC) curves of the*

309     *individual 3D reconstructions computed from particles labeled "good" (also see **c**). **(e)***

310     *Representative alpha-helix (amino acids 600-615) illustrating the improvement of the*

311     *density when using the final (**bottom**) compared to the initial (**top**) picking model. **(f)***

312     *3D reconstruction of TRPC4 computed from 500 micrographs using the optimized*

313     *picking model.*

314

315     the first iteration of the feedback loop (Figure 3a). This was done by replacing 90% of

316     the particle boxes determined by crYOLO with randomly positioned boxes within the

317     same micrographs. In combination, these measures ensured that the initial picking

318     results were almost entirely unusable and successful re-training had to take place in

319     order to enable further processing of the data.

320          Despite the bad starting point, by the final feedback loop iteration the

321     repeatedly re-trained model has successfully learned to pick the previously unknown

322     TRPC4 particles resulting in high-resolution 2D class averages (Figure 3a, b). An

323     evaluation of the performance of the feedback loop on a fixed subset of 500

324     micrographs (see Methods for details), illustrates that the number of "good" particles

325     increases sharply within the early iterations of the feedback loop from an initial 25%

326     of particles to a stable value of ~ 50%, (Table 1, Figure 3c) and a final resolution of 3.6

327     Å (FSC=0.143). This increased ability to identify a greater number of usable particles

328     on the same subset of micrographs is also reflected in the map quality and achieved

329     resolution when using the intermediate crYOLO models produced during the

330     individual feedback rounds to process the fixed set of 500 micrographs (Figure 3d-f).

331          This experiment furthermore demonstrates the ability of crYOLO to adapt to

332     unknown data even if only sparse training data is available. In the initial round of the

333     feedback loop a mere 5 particles per micrograph ended up in "good" classes on

334    average – and, consequently, are all that was available to re-train the picking model

335    (Table 1).

| Feedback round | Good classes | Good particles | Picks/Mic | Good picks/Mic | Resolution | Relative good picks |
|---|---|---|---|---|---|---|
| 2 | 28 | 2,504 | 20 | 5 | 4.71 | 0.25 |
| 3 | 236 | 22,154 | 104 | 44 | 3.65 | 0.43 |
| 4 | 349 | 32,513 | 132 | 65 | 3.60 | 0.49 |
| 5 | 355 | 33,415 | 152 | 67 | 3.60 | 0.44 |
| 6 | 361 | 34,018 | 147 | 68 | 3.55 | 0.46 |
| 6 + T 0.257 | 331 | 31,237 | 114 | 62 | 3.6 | 0.55 |

336

337    ***Table 1. TRPC4 feedback loop statistics.*** *For every feedback round as well as the final*

338    *run after optimization of the picking threshold (6 + T x.xx) the number of classes*

339    *labeled "good" by Cinderella; the number of particles included in these classes; the*

340    *total number and the number of good particles picked per micrograph; the final*

341    *resolution of the 3D reconstruction; and the relative amount of good particles are*

342    *listed for the TRPC4 data (500 micrographs).*

343

344    In summary, the TranSPHIRE feedback loop is able to automatically optimize

345    the internally used picking model and provide reliable, high quality picking results

346    even when processing challenging samples that initially are barely recognized by the

347    model. We have shown that in such a case, after five feedback rounds, crYOLO is able

348    to pick the TRPC4 membrane protein to completion, without requiring the user to

349    continuously monitor, let alone disrupt the ongoing data processing. The feedback

350    loop optimization is fully integrated into the TranSPHIRE pipeline and works entirely

351    automated out of the box. Its capabilities extend to difficult data sets such as

352    membrane proteins, and enable advanced processing methods, such as targeting

353    specific conformational states, or processing filamentous data sets, as demonstrated

354    in the following.

355

356

357 **Selectively targeting a conformational state in a mixed sample using prior**

358 **knowledge**

359 A basic assumption of most algorithms currently used to process cryo-EM data is that

360 all particles in a data set are projections of the same structure, hidden behind a

361 curtain of noise. In reality, however, cryo-EM samples are often more complex, and

362 can contain multiple conformational states of the target structure, impurities, and

363 aggregates. Filtering such unwanted data and selectively targeting only a subset of

364 the structures found within a sample is one of the fundamental issues in cryo-EM,

365 and often requires significant efforts to address and resolve.

366 The TranSPHIRE feedback loop offers a straightforward solution to this issue by

367 allowing the injection of additional knowledge into the pipeline, either before or

368 during runtime. This enables users to incorporate and make use of information that is

369 already available, as well as information that was just produced during acquisition.

370 Specifically, a set of 2D class averages of the target structure can be used to train

371 Cinderella [20] to only recognize these averages as representatives of "good" classes,

372 and, consequently, everything else as "bad." If such averages are available

373 beforehand, Cinderella can be pre-trained; otherwise the feedback loop can be

374 paused once the first set of 2D class averages are produced in the TranSPHIRE

375 pipeline and continued after manual re-training of Cinderella. This additional training

376 step to embed additional knowledge into the TranSPHIRE pipeline enables us to steer

377 the re-training of the picking model during the feedback loop iterations. More

378 precisely, particles that end up in sharp classes depicting a different particle, a

379 subcomplex, and/or the target protein in the wrong conformational state (for

380 example) will now also be labeled as "bad" by Cinderella, despite their high quality.

381 During the feedback loop, crYOLO will thus be taught to only focus on particles that

16

382  end up in quality classes depicting the wanted particle or state, while, at the same

383  time, reject anything else, including sharp classes from an unwanted subpopulation.

384  To demonstrate the capability of the TranSPHIRE feedback loop to use prior

385  knowledge and target a pre-selected conformation, we processed a sample of the Tc

386  holotoxin that contained particles in two conformational states, namely the pre-pore

387  and pore state (Figure 4a). Of these, we only targeted the pore state, which is

388  significantly more difficult to find as it only accounts for ~ 19% of the particles within

389  the data set (Table 2, Figure 4b). Cinderella was trained with 318 examples of "good"

390  classes (side-views of the pore state) and 664 examples of "bad" classes (views of the

391  pre-pore state and contamination). During the feedback loop crYOLO was then re-

392  trained with only those particle picks that ended up in "good" classes showing views

393  of the pore state.

| Feedback round | Good classes | Good particles | Picks/Mic | Good picks/Mic | Resolution | Relative good picks |
|---|---|---|---|---|---|---|
| 1 | 130 | 12,595 | 130 | 25 | 4.28 | **0.19** |
| 2 | 145 | 10,406 | 100 | 21 | 4.24 | **0.21** |
| 3 | 151 | 14,534 | 74 | 29 | 4.36 | **0.40** |
| 4 | 146 | 14,081 | 71 | 28 | 4.28 | **0.40** |
| 5 | 155 | 14,935 | 68 | 27 | 4.28 | **0.40** |
| 6 | 140 | 13,566 | 69 | 27 | 4.24 | **0.39** |
| 6 + T 0.194 | 145 | 13,954 | 55 | 28 | 4.24 | **0.50** |

394

395  ***Table 2. Tc holotoxin feedback loop statistics.*** *For every feedback round as well as*

396  *the final run after optimization of the picking threshold (6 + T x.xx) the number of*

397  *classes labeled "good" by Cinderella; the number of particles included in these*

398  *classes; the total number and the number of good particles picked per micrograph;*

399  *the final resolution of the 3D reconstruction; and the relative amount of good*

400  *particles are listed for the Tc Holotoxin data (500 micrographs).*

401

402          To evaluate the performance of the feedback loop we used the intermediate

403  picking models produced during the individual feedback rounds to separately process

404  a fixed set of 500 micrographs once the feedback loop had finished (see Methods for

17

405   details). We observed a steady decrease of particles representing the pre-pore state

406   – that we are not interested in – together with an initial rise and then level amount of

407   pore state picks (Figure 4b, Table 2). While initially only 19% of the particles

408   resembled the pore state, slightly more than 50% of all picks ended up in 2D class

409   averages depicting our targeted conformation when using the final optimized picking

410   model (Figure 4c-d). As in the previous experiment, the percentage of relative good

411   picks per micrograph steadily increases. Notably, this happens while neither the

412   number of good classes, nor the number of good particles seem to follow suit (Table

413   2). This means that our re-training efforts are working as intended: Over the course

414   of the feedback loop, crYOLO learns to discard quality class averages of the pre-pore

415   state that we are not interested in and instead focus on picking the less common

416   pore state. Consequently, the amount of picked particles changes slowly, while, at

417   the same time, the relative amount of "good" particle picks steadily increases,

418   resulting in a 4.2 Å (FSC=0.143) 3D reconstruction of the pore state from no more

419   than 500 micrographs (Figure 4e).

420       Taken together, these results illustrate how additional knowledge can be

421   used to pre-train Cinderella, allowing TranSPHIRE to steer the re-training of the

422   picking model during the feedback loop and to target a known subpopulation within

423   the data. Using a picking model optimized for a specific conformation offers a two-

424   fold advantage. First, reconstruction efforts will be more effective, as we gain more

425   particles of the subpopulation that we are interested in. Second, reconstruction

426   efforts will be more efficient, as the rejection of particles that end up in "bad" classes

427   significantly shrinks the overall size of the data set. In our example, we reduce the

428   number of picked particles from an initial total of 67,117 to a set of only 27,646

429   particles, without reducing the achieved resolution or the number of pore state

430   particles that we are interested in (Figure 4b). Any follow-up computations, such as

431    costly 3D reconstructions, benefit greatly from such a reduction in data set size as it

432    results in a much more efficient use of the available computational resources.
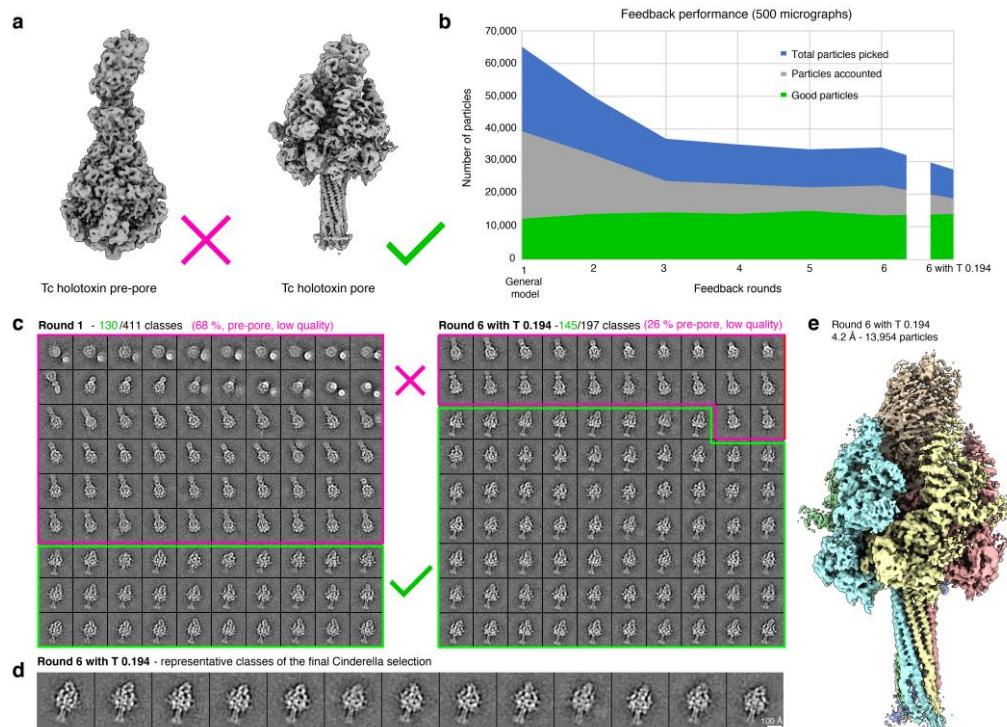
433



434

435    ***Figure 4. Using prior knowledge to extract a pre-selected conformational state. (a)***

436    *The processed data set contains the Tc holotoxin in both the pre-pore state (**left**) and*

437    *the more rare pore state (**right**). In this experiment we specifically target the pore*

438    *state. **(b)** Progression of the number of picked particles (blue), those accounted during*

439    *2D classification (grey) and particles labeled "good" i.e. representing the pore state*

440    *(green) when applying the intermediate picking models of the feedback loop to a*

441    *fixed subset of 500 micrographs. Initial picking is dominated by pre-pore state*

442    *particles. This overhead is reduced with each iteration, while the amount of picked*

443    *pore state particle remains stable. **(c)** Representative 2D class averages depicting the*

444    *decrease of unwanted classes (pore state or low quality; marked red) from an initial*

445    *68% in the first feedback round (**left**) to 26% after the last feedback round (**right**). **(d)***

446    *Representative 2D class averages depicting the pore state as selected by Cinderella in*

447    *the final iteration of the feedback loop. **(e)** 3D reconstruction of the Tc holotoxin pore*

448    *state computed from 500 micrographs using the final optimized picking model.*

449

450

19

451 **Using TranSPHIRE to automatically process filamentous proteins**

452 Filamentous proteins such as the actomyosin complex are notoriously difficult to

453 process. This is because their structure is by definition not limited to a single element

454 but rather forms a continuous strand that both enters and exits the enclosing frame

455 of any picked particle image. Consequently, filamentous proteins are traced, rather

456 than picked, and overlapping segments have to be identified along each filament,

457 while filament crossings and contamination need to be avoided. In addition,

458 filamentous projections share a similar overall geometry which increases the

459 correlation between any two particles and interferes with alignment attempts during

460 2D classification. While there are several programs available that implement manual

461 filament processing [13,24-27], until now there has not yet been any cryo-EM software

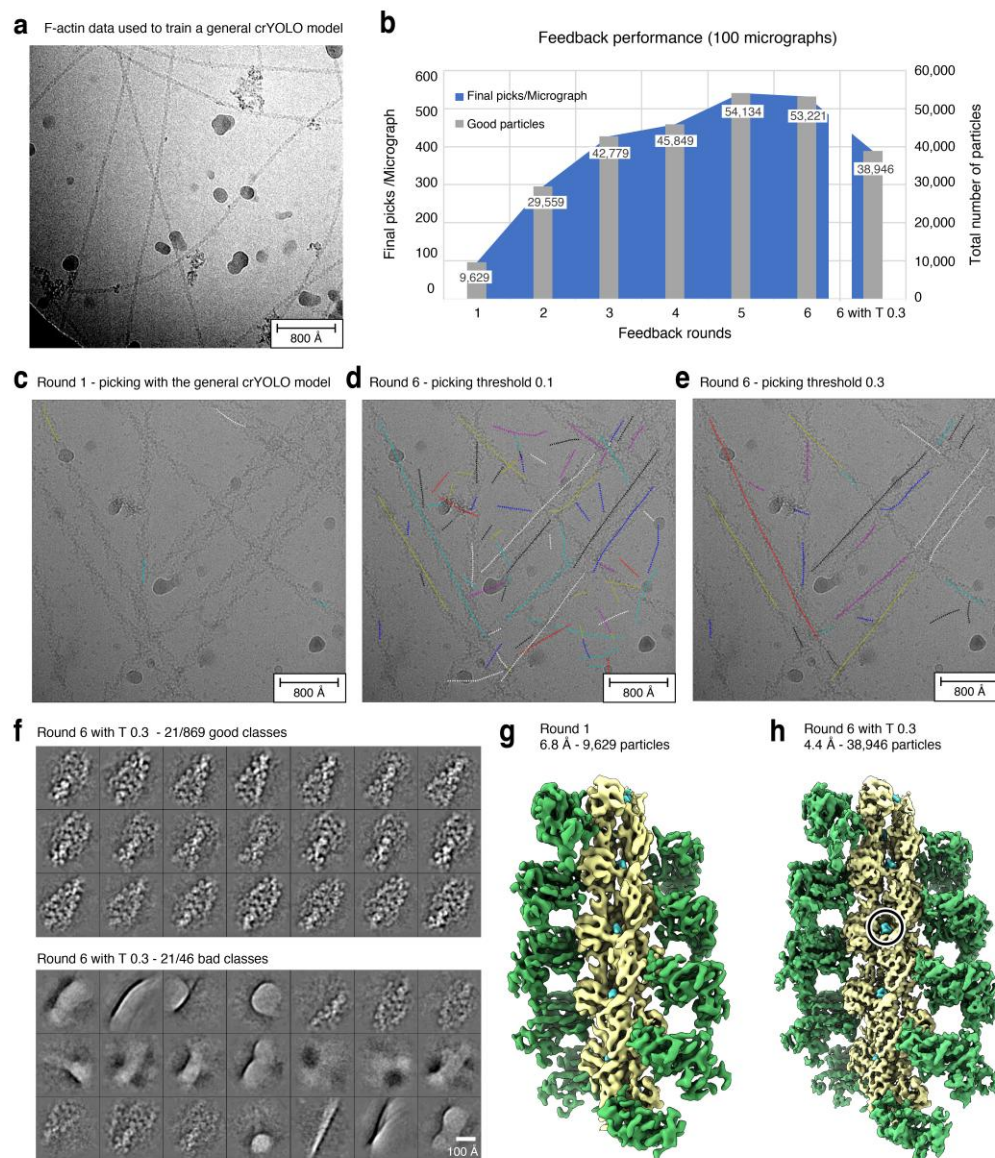462 package that offers the automated processing of filamentous data sets.

463 With TranSPHIRE we introduce a comprehensive software package for cryo-

464 EM that includes the ability to automatically process filamentous proteins utilizing

465 methods of the SPHIRE package [13]. While the actual processing is fully automated,

466 some preparation is still needed when using the TranSPHIRE pipeline to process

467 filaments. Specifically, crYOLO needs to be trained to pick filaments [28], as these look

468 fundamentally different from the single particle complexes known to its default

469 general model. Additionally, Cinderella [20] also needs to be trained with 2D class

470 averages of the filament in question. If such class averages are not available initially,

471 the feedback loop can be halted for re-training Cinderella as soon as TranSPHIRE has

472 produced them. Once the models for the deep learning decision makers of the

473 pipeline are trained on the specific filamentous data, TranSPHIRE and its integrated

474 feedback loop are ready to automatically process the respective filamentous data

475 sets.

476        As an example of processing initially unknown filamentous data, we chose an

477    actomyosin complex. To further demonstrate the ability of the feedback loop to

478    adjust the picking to a specific filamentous protein complex, we trained crYOLO with

479    multiple data sets of F-actin, which looks substantially different than the actomyosin

480    complex (Figure 5a). Thereby, crYOLO learns to trace filaments, but does not readily

481    recognize actomyosin filaments resulting in a weak initial picking performance

482    (Figure 5b-c).

483        As soon as the first 2D class averages became available, the feedback loop

484    was halted and a new Cinderella model was trained manually. Afterwards the

485    feedback loop continued through its default five iterations, automatically teaching

486    crYOLO to identify projections of the actomyosin complex. To evaluate the

487    performance of the feedback loop we separately processed a fixed set of 100

488    micrographs using the intermediate picking models produced during the individual

489    feedback iterations (Figure 5, Table 3, see Methods for details).

| Feedback round | Good classes | Good particles | Picks/Mic | Good picks/Mic | Resolution | Relative good picks |
|---|---|---|---|---|---|---|
| 1 | 211 | 9,629 | 143 | 96 | 7.63 | **0.67** |
| 2 | 659 | 29,559 | 358 | 296 | 4.48 | **0.83** |
| 3 | 936 | 42,779 | 552 | 428 | 4.37 | **0.77** |
| 4 | 1,016 | 45,849 | 639 | 458 | 4.54 | **0.72** |
| 5 | 1,203 | 54,134 | 1,098 | 541 | 4.32 | **0.49** |
| 6 | 1,174 | 53,221 | 1,073 | 532 | 4.72 | **0.50** |
| 6 + T 0.3 | 869 | 38,946 | 515 | 389 | 4.54 | **0.76** |

490

491    ***Table 3. Actomyosin complex feedback loop statistics.*** *For every feedback round as*

492    *well as the final run after optimization of the picking threshold (6 + T x.xx) the number*

493    *of classes labeled "good" by Cinderella; the number of particles included in these*

494    *classes; the total number and the number of good particles picked per micrograph;*

495    *the final resolution of the 3D reconstruction; and the relative amount of good*

496    *particles are listed for the actomyosin complex data (100 micrographs).*

497

498

**Figure 5. Ligand identification within an actomyosin complex. (a)** *Representative micrograph of the F-actin data used to train crYOLO. **(b)** Progression of the number of "good" particles per micrograph (blue) and in total (grey) when applying the intermediate picking models of the feedback loop to a fixed subset of 100 micrographs. The dipping curve at the end indicates the desired loss of low-quality picks that are excluded when a higher picking threshold (0.3) is used. **(c)** Representative micrograph of the actomyosin complex highlighting the weak initial picking results when using the crYOLO model trained on F-actin data (see **a**). **(d)** Particle picking performance on the same micrograph using the final picking model. While filaments are now traced much more effectively, the model also picks unwanted filament crossings and contamination. **(e)** Increasing the picking threshold from 0.1 to the default value of 0.3 minimizes the amount of false positive picks, while*

22

511   *maintaining the desired filament traces. **(f)** Representative 2D class averages labeled*

512   *"good" **(top)** and "bad" **(bottom)** by Cinderella based on 100 micrographs and using*

513   *the final model for picking. **(g)** 3D reconstruction of the actomyosin complex*

514   *computed from 100 micrographs using the initial picking model. **(h)** 3D reconstruction*

515   *computed from the same 100 micrographs using the final optimized picking model.*

516   *The resolution is sufficient to verify the binding of a ligand (circled).*

517

518   Initially a low confidence threshold of 0.1 (default) was used for picking in

519   order to gather enough training data (Figure 5c, d). However, the amount of picked

520   particles and the confidence in the picks increased throughout the feedback loop

521   (Figure 5b, c). Thus, the picking threshold was adjusted to the default value of 0.3

522   after the feedback in order to exclude low confidence picks of contamination and

523   filament crossings (Figure 5b, d-e). Thereby, the number of relative good particles

524   could be increased from 50% to 76% (Table 3) resulting in few classes labeled "bad"

525   (Figure 5f). The improvement is also visible when comparing the initial and final 3D

526   reconstruction computed from the same set of 100 micrographs (Figure 5g-h).

527   Particularly, the final reconstruction of 4.4 Å (FSC=0.143) is sufficient to identify a

528   small molecule bound to the filament, highlighting how TranSPHIRE can simplify

529   ligand screenings.

530   Using the feedback loop, TranSPHIRE offers the first cryo-EM software

531   package that is able to automatically process filamentous data, even if the precise

532   shape of a specific filament is initially unknown to the pipeline. Moreover,

533   TranSPHIRE now enables experimentalists to produce an early 3D reconstruction with

534   a resolution sufficient to identify bound ligands and determine whether or not their

535   data is likely to yield a high resolution reconstruction – all within the time frame of

536   hours and while their data is still being collected at the microscope (Figure 2,

537   Supplementary Figure 3). The automated processing greatly simplifies the processing

23

538    of filamentous samples in general and, most importantly, facilitates the fast

539    determination of multiple structures of one filament decorated with different

540    accessory proteins or bound to ligands.

541

## 542    **4     Discussion**

543    In this paper we present the streamlined TranSPHIRE pipeline for automated,

544    feedback-driven processing of cryo-EM data. It fully automates data transfer, pre-

545    processing and the creation of a series of early reconstructions based on the most

546    recently processed data (Figure 1a). At the same time, TranSPHIRE prominently

547    displays all relevant data evaluation metrics, updated in real time (Supplementary

548    Figure 1), and offers the option to send email notifications when issues are

549    encountered or important milestones – such as the first 2D class averages, or an

550    initial 3D reconstruction – have been reached.

551        We also introduce the TranSPHIRE feedback loop (Figure 1b), a machine

552    learning-based method to optimize the internally used particle picking model and

553    adapt our native crYOLO picker to any data set, even while it is still being collected at

554    the microscope. This allows TranSPHIRE to adjust to never before seen data, as well

555    as to avoid any issues that a cryo-EM sample might include, such as unwanted

556    proteins, low contrast, and/or different kinds of contamination. The optimization of

557    the picking model performed by the feedback loop can further be guided by the

558    experimentalist in order to specifically select a subpopulation within the data, such as

559    a distinct conformational or oligomeric state.

560        We demonstrate these capabilities of TranSPHIRE and its new feedback loop

561    by performing three distinct experiments, each addressing a common issue in cryo-

562    EM: First, we processed the membrane protein TRPC4, after purposefully sabotaging

563    our particle picking to simulate processing a data set that is not only unknown, but

24

564    initially only barely provides enough useful picks for training. Nevertheless, the

565    TranSPHIRE feedback loop successfully taught crYOLO to identify and pick the sought-

566    after particles without any need for user intervention or expert knowledge input.

567    When the final picking model was then used to automatically compute a full

568    reconstruction, we reached a resolution of 3.6 Å (FSC=0.143), based on the data

569    extracted from no more than 500 micrographs (Figure 3).

570        Second, we processed a sample containing the Tc holotoxin in two different

571    conformational states: The common pre-pore state, and the significantly rarer pore

572    state that only accounts for about one fifth of the available particles. In this

573    experiment we injected prior knowledge about the pore state into the pipeline by

574    training Cinderella – our deep learning tool to reject unusable 2D class averages – to

575    only accept class averages of this state. This directed the re-training during the

576    feedback loop and taught crYOLO to focus on the rare pore state particles. As a

577    result, we obtained a picking model that was highly selective for only one

578    conformational state while rejecting not only low quality 2D class averages, but also

579    high quality 2D class averages if they displayed the Tc holotoxin in the

580    conformational state that we were not interested in (Figure 4). This produced a

581    particle stack that was not only populated with an increased number of "good"

582    particles, but also contained less particles overall, as unwanted particles were already

583    rejected during particle picking. The final reconstruction obtained a resolution of 4.2

584    Å (FSC=0.143). Such an optimized stack means that any follow-up computations only

585    have to deal with relevant data, allowing for a more efficient use of the available

586    computational resources.

587        Third, we processed a data set of an actomyosin complex to demonstrate

588    how the ability of TranSPHIRE to automatically process cryo-EM data also extends to

589    filamentous proteins. To adjust the pipeline to the processing of filaments, we re-

590  trained both crYOLO and Cinderella in order to teach them about the distinct visual

591  properties of filamentous particles and how to avoid any filament-exclusive pitfalls,

592  such as filament crossings. To specifically showcase the ability of the feedback loop

593  to deal with an initially unknown filament structure, we only taught crYOLO about F-

594  actin, which features a fundamentally different appearance than the actomyosin

595  complex. Cinderella was then only trained with the initial 2D class averages that

596  TranSPHIRE produced during the first iteration of the feedback loop. Despite the

597  initial picking model only knowing about F-actin, Cinderella was able to teach crYOLO

598  about the actomyosin complex and the final reconstruction reached a resolution of

599  4.4 Å (FSC=0.143), using the data extracted from merely 100 micrographs (Figure 5).

600  In summary, TranSPHIRE offers a fully automated pipeline that produces

601  highly optimized particle stacks that allow for more effective processing and more

602  efficient use of any available resources, both computational and human. Combined,

603  these features allow experimentalists to make the most of their limited time at the

604  microscope and to identify and address any issues as soon as they surface.

605  Furthermore, TranSPHIRE produces early reconstructions of proteins, even if initially

606  unknown, thereby enabling experimentalists to assess their data and identify the

607  conformational state of their protein or validate the binding of a ligand while their

608  data is still being collected. Hence, TranSPHIRE allows users to perform automated

609  high-throughput on-the-fly screenings for different buffer conditions or ligands of

610  interest.

611  

## Methods

612  

613  **Hardware used to run TranSPHIRE**

614  By default, TranSPHIRE runs on a single machine, which can be combined with a

615  separate workstation or computer cluster to outsource computational power. For the

616    majority of the results presented in this manuscript, a single machine equipped with

617    two Intel(R) Xeon(R) Gold 6128 CPUs (3.40GHz), featuring 12 CPU cores each

618    (hyperthreading 24); 192 GB of RAM; and three GeForce RTX 1080 Ti GPUs was used.

619    Only computationally more expensive 3D reconstructions, both the initial *ab initio*

620    reconstruction and the 3D refinements (for details see below), were outsourced to

621    our local computer cluster. There, calculations were performed on two nodes; each

622    equipped with two Intel(R) Xeon(R) Gold 6134 CPUs (3.20GHz), featuring 32 CPU

623    cores in total and 384 GB of RAM.

624

625    **Software integrated into the TranSPHIRE pipeline**

626    TranSPHIRE is a free of charge, open-source software written in Python3, which is

627    available online (https://github.com/MPI-Dortmund/transphire).

628      Its fully-automated processing pipeline integrates several software packages

629    and is thereby highly flexible and adaptable. An initial integrity check and the

630    consecutive compression of every input stack to a LZW compressed tiff file is

631    performed using IMOD v4.9.8 [30]. Currently, TranSPHIRE supports several options for

632    motion correction (Unblur [15] and MotionCor2 [14]) and CTF estimation (CTFFIND [16],

633    CTER [17] and GCTF [18]). For all consecutive 2D and 3D processing steps, TranSPHIRE

634    utilizes functions of the SPHIRE [13] package including the deep-learning particle picker

635    crYOLO [12], the 2D class selection tool Cinderella [20] and a new GPU accelerated version

636    of the reliable 2D classifier ISAC2 [19].

637      Results presented in this manuscript were generated with TranSPHIRE

638    v1.4.50 and SPHIRE v1.4. Specifically, the pipeline consisted of the following modules:

639    the CUDA 10.2.86 version of MotionCor2 v1.3.0[14]; CTFFIND v4.1.13 for CTF

640    estimation [16]; crYOLO v1.6 for particle picking [12]; SPHIRE sp_window.py for particle

641    extraction [13]; a GPU accelerated version (v1.0) of SPHIRE ISAC2 [19] for on-the-fly 2D

642     classification (will be published elsewhere); SPHIRE Cinderella v0.5 [20] for 2D class

643     selection; SPHIRE sp_rviper.py [13,21] for *ab initio* reconstructions and finally SPHIRE

644     sp_meridien.py [13] or sp_meridien_alpha.py for the 3D refinement of single particles

645     or filaments, respectively.

646

647     **The automated processing pipeline within TranSPHIRE**

648     After preprocessing the data i.e. data transfer and compression, motion correction

649     and CTF estimation (also see Supplementary Figure 2), particles are automatically

650     picked using the deep learning, GPU-accelerated particle picker crYOLO [12]. By using

651     the general model, which was trained on 63 cryo-EM data sets, crYOLO is able to pick

652     previously unseen particles. During the feedback rounds a picking threshold of 0.1 is

653     used to facilitate the picking of distinct proteins and features. At the end of each

654     feedback iteration crYOLO is retrained on particles that contributed to classes labeled

655     "good" by Cinderella (see below and Figure 1b). When crYOLO is trained on a single

656     data set, it quickly reaches a good picking quality even when the training data only

657     contains few micrographs. Hence, increasing the size of the training data, enhances

658     the training time without benefitting the training. Therefore, only particles from 50

659     randomly selected micrographs and no more than 20,000 particles in total are used

660     for the training. Once the feedback loop is finalized, the picking performance is

661     further optimized by adjusting the picking threshold to an optimal one, as

662     determined by a parameter grid search using crYOLO's internal evaluation procedure.

663     The particle threshold value defines a confidence threshold that each pick made by

664     crYOLO must either meet or exceed in order to be accepted. If this threshold is set to

665     a low value, particles with a low confidence are also accepted. In order to find the

666     optimal threshold, a fixed subset of data is repeatedly picked while varying the

667     threshold from 0.0 to 1.0, using a step size of 0.01. Afterwards the optimal threshold

668    is defined by the highest F2 score [31] of all resulting picks. Processing results

669    generated with the optimized threshold are labeled with iteration "6 + T x.xxx",

670    where six represents the sixth and thus final model used in the feedback loop, and

671    the value x.xxx denotes the optimized picking threshold.

672         Picked particles are automatically extracted and classified in 2D, resulting in

673    class averages containing 60 to 100 particles per class (standard settings).

674    Classifications are performed by a new GPU-accelerated and updated version of

675    ISAC2, which is based on the original ISAC (Iterative Stable Alignment and Clustering)

676    algorithm[19]. Just like the CPU-bound ISAC2 it delivers high quality 2D class averages

677    as well as an initial clean-up of the data set, but does not come with the same high

678    computational cost. Hence, GPU ISAC provides the same functionality on a single

679    workstation without the need to outsource 2D classification to a cluster. The GPU

680    ISAC code repository is part of the SPHIRE repository listed above.

681     As the generation of high-resolution 2D class averages requires a sufficient number

682    of particles covering a range of views, 2D classification is only started once a certain

683    number of particles is accumulated. While this number can be adjusted in the

684    TranSPHIRE GUI, a default value of 20,000 particles per batch has proven to be good

685    (see also Supplementary Figure 4).

686         2D class averages are routinely used to assess the overall quality of the data

687    and to select only those particles for 3D refinement that contribute to high quality 2D

688    class averages. Previously, this selection was done manually, breaking any automated

689    processing pipeline. In order to provide a fully automated pipeline, TranSPHIRE uses

690    Cinderella [20], a deep learning binary classifier based on a convolutional neural

691    network. When provided with a set of 2D class averages, Cinderella labels each of

692    them as either "good" or "bad." By default, this decision is based on a model that was

693    trained on a large set of class averages from a multitude of different cryo-EM

694    projects. Alternatively, Cinderella can be trained on specific data to select classes

695    according to the needs of the current project. By default, TranSPHIRE runs Cinderella

696    using its general model, based on 3,559 "good" and 2,433 "bad" classes taken from

697    20 different data sets from both the EMPIAR [32] data base and our in-house efforts.

698    The Cinderella git repository can be found online

699    (https://github.com/MPI- Dortmund/sphire_classes_autoselect).

700    Once the feedback loop has finished and a set of at least 200 "good" class

701    averages is available (number can be adjusted if desired), a reproducible, *ab initio* 3D

702    reconstruction is computed from 2D class averages using the SPHIRE method RVIPER

703    [13] (Reproducible Validation of Individual Parameter Reproducibility). The VIPER

704    algorithm combines a genetic algorithm [33] with stochastic hill climbing [34] to produce

705    multiple 3D ab initio structures. These reconstructions are then compared and the

706    most reproducible model is used to seed the consecutive 3D refinement. (See online

707    documentation for RVIPER and VIPER at

708    http://sphire.mpg.de/wiki/doku.php?id=pipeline:viper:sxrviper).

709    To generate a high-resolution 3D reconstruction a stack of all particles

710    assigned to classes that were labeled "good" by Cinderella is created. The

711    consecutive refinement is performed by the SPHIRE method MERIDIEN [13] providing

712    the initial reconstruction computed in the previous step as reference. The refinement

713    within MERIDIEN proceeds in two phases. The first phase, "EXHAUSTIVE", searches

714    the whole 3D parameter space -- three Euler angles for rotation and two dimensions

715    for translation -- on a discrete grid. The second phase, "RESTRICTED", searches the

716    parameter space on a discrete grid within the local area closest to the best matching

717    set of parameters found in the previous iteration. To avoid over-fitting, the image

718    dimensions and the grid spacing is adjusted after every iteration, based on the

719    achieved resolution according to the gold standard FSC [35] and stability of the

720    parameters. In order to compensate for the discreteness of the grid and the

721    uncertainty in parameter assignment, particles are weighted by the probability of the

722    parameter set for the backprojection into the 3D reconstruction. (See online

723    documentation of MERIDIEN at

724    http://sphire.mpg.de/wiki/doku.php?id=pipeline:meridien:sxmeridien).

725        Similar to the prerequisites for 2D classification, a certain number of particles

726    representing different views is required to successfully compute a 3D reconstruction.

727    Thus, TranSPHIRE will not start the 3D refinement before a defined number of

728    particles is accumulated. In our hands a total of 40,000 particles (default value, can

729    be adjusted) is sufficient to calculate a medium to high resolution 3D reconstruction

730    in a short time frame. While this reconstruction will likely not reach the highest

731    resolution possible, it still enables a first analysis i.e. identification of a

732    conformational state or the verification if a ligand is bound or not. Furthermore, it

733    provides a quality control throughout the data acquisition, as a new 3D

734    reconstruction is computed for every batch of 40,000 particles. As all 3D refinements

735    start from the same initial reference, refinement projections parameters can

736    additionally be used to directly start with a local refinement of the complete data set,

737    thereby significantly reducing the required running time.

738

739    **Evaluation of the feedback performance**

740    As TranSPHIRE runs in parallel to the data acquisition and data are processed as they

741    come in, the number of movies is increasing during the runtime and results from one

742    feedback iteration to the next are not directly comparable. Thus, the feedback

743    performance was evaluated separately for every data set on a fixed subset of 500

744    (TRPC4 and Tc holotoxin, Figure 3-4) and 100 (Actomyosin, Figure 5) micrographs.

745   For each case, the fixed subset was processed using the intermediate picking models

746   produced during the individual feedback iterations. Specifically, every subset was

747   once picked with the starting model (general model, labeled round 1) and with every

748   picking model generated throughout the five iterations of the feedback loop (rounds

749   2 to 6) using a particle threshold of 0.1. In addition, another run was performed with

750   the final picking model using the optimized particle threshold (6 + T X.XX). The

751   consecutive processing in 2D and 3D was performed with AutoSPHIRE sp_auto.py,

752   which is the automatic, batch processing  tool within SPHIRE [13] on our local CPU

753   cluster. The processing pipeline and settings used resemble the ones described

754   above, except that CPU ISAC was used instead of the new GPU-accelerated version.

755

756

757   **Automatic processing of the TRPC4 data.**

758   The performance of TranSPHIRE was tested on a subset of 500 micrographs of a high-

759   resolution data set of the transient receptor channel 4 (TRPC4) from zebra fish in

760   LMNG detergent (prepared in analogy to [36], publication in preparation). The data set

761   was automatically collected at a Cs-corrected Titan Krios (FEI Thermo Fisher)

762   microscope equipped with an X-FEG and operated at 300kV using EPU (FEI Thermo

763   Fisher). Equally dosed frames with a pixel size of 0.85 Å/pixel were collected with a

764   K2 Summit (counting mode, Gatan) direct electron detector in combination with a

765   GIF quantum-energy filter set to a filter width of 20 eV.  Each movie contains 50

766   frames and a total electron dose of 88.5 e/$Å^2$.

767        Processing in TranSPHIRE was performed as described above with five

768   internal feedback rounds to optimize the crYOLO picking model. Within the pipeline,

769   movies were drift corrected and dose weighted by MotionCor2 [14] using five patches

770   with an overlap of 20% and CTFFIND4 [16] fitted the CTF between 4 Å and 30 Å with an

771  Cs value of 0.001. The training data for the general model of crYOLO usually contain

772  four data sets of TRP channels. To avoid any favorable picking bias and handle the

773  TRPC4 data as previously unseen, the general model was retrained after removing all

774  TRP channels from the training data.  Even then, crYOLO was able to identify most

775  TRPC4 particles through the successful generalization. To simulate a worst-case

776  scenario of a deficient initial picking performance, 90% of the particle picks in the

777  initial feedback round were replaced by random coordinates.

778  During the feedback rounds the crYOLO picking threshold was set to 0.1 and

779  the anchor size to the estimated particle diameter of 240 pixels. After the final

780  feedback round, the picking threshold value was adjusted to 0.257 based on the

781  crYOLO confidence threshold optimizing procedure described above. After each

782  particle picking step, particles were automatically extracted using SPHIRE

783  sp_window.py with a box size of 288 pixels. The subsequent 2D classification was

784  performed using a GPU accelerated version of the SPHIRE ISAC2 algorithm using

785  standard settings. The feedback loop was run with the default particle batch size of

786  20,000 (for details see above and Supplementary Figure 3).

787  The produced 2D class averages were subjected to an automatic 2D class

788  selection using our deep learning tool Cinderella and a confidence threshold of 0.1.

789  To simulate the processing of a previously unseen protein, Cinderella was trained

790  with its general model training data excluding all channel proteins, thereby ensuring

791  an unbiased selection process. During the feedback rounds crYOLO was trained on

792  the default value of 50 random micrographs that contained particles contributing to

793  classes labeled "good" by Cinderella. 3D reconstructions were computed as described

794  above using no mask and imposing c4 symmetry. Note that albeit our program

795  provides the possibility to compute a 3D mask from the initial model automatically

796  and apply it during the refinement, this option is deactivated by default. Automated

797    masking procedures might eliminate valid regions of the structure that are not well

798    resolved in the initial reconstruction, especially in cases with strong flexibility in the

799    complex. In case a 3D mask is not provided, we strongly recommend to use a mask

800    created from the results of TranSPHIRE for all follow-up experiments, in order to

801    exploit the full potential of 3D refinement. Whereas the workflow can be easily

802    extended, the pipeline for each batch stops by default after the first high resolution

803    3D refinement, in order to allow on-the-fly evaluation by the user. The results can be

804    easily converted to RELION after any milestone and *vice versa*. Correction of higher-

805    order aberrations for example in RELION  might further improve the resolution of the

806    final result, when these optical effects  are present [37].

807        The progression of the picking performance throughout the feedback rounds

808    was evaluated on a fixed subset of 500 micrographs as described above (Figure 3).

809    Note that the picking model of the first iteration is not included in this evaluation, as

810    its performance was initially corrupted by randomizing 90% of the picked particles.

811

812    **Automatic processing of the Tc holotoxin data.**

813    To test the capability of TranSPHIRE to target a specific conformation, a subset of 500

814    micrographs of the ABC holotoxin from *Photorhabdus Luminescens* reconstituted in a

815    lipid nanodisc (EMD-10313) [29] was processed. This data set contains a mixture of

816    conformations, namely the pre-pore and pore state of the holotoxin. The data set

817    was collected at a Cs-corrected Titan Krios (FEI Thermo Fisher) microscope equipped

818    with an X-FEG and operated at 300kV using EPU (FEI Thermo Fisher). Equally dosed

819    frames with a pixel size of 0.525 Å/pixel were collected with a K2 Summit (super

820    resolution mode, Gatan) direct electron detector in combination with a GIF quantum-

821    energy filter set to a filter width of 20 eV.  Each movie contains 40 frames and a total

822    electron dose of 60.8 e/Å$^2$.

823       Processing in TranSPHIRE was performed as described above with five

824    internal feedback rounds to optimize the crYOLO picking model. Within the pipeline,

825    movies were drift corrected, dose weighted and binned to a pixel size of 1.05 Å/px by

826    MotionCor2 [14] using three patches without overlap and CTFFIND4 [16] fitted the CTF

827    between 4 Å and 30 Å with an Cs value of 0.001. Subsequently, particles were picked

828    using the general model of crYOLO.

829    During the feedback rounds the crYOLO picking threshold was set to 0.1 and the

830    anchor size to the estimated particle diameter of 205 pixels. After the final feedback

831    round, the picking threshold value was adjusted to 0.194 based on the crYOLO

832    confidence threshold optimizing procedure described above. After each particle

833    picking step, particles were automatically extracted using SPHIRE sp_window.py with

834    a box size of 420 pixels. The subsequent 2D classification was performed using a GPU

835    accelerated version of the SPHIRE ISAC2 algorithm using standard settings. The

836    feedback loop was run with the default particle batch size of 20,000 (for details see

837    above and Supplementary Figure 3).

838       The produced 2D class averages were subjected to an automatic 2D class

839    selection using our deep learning tool Cinderella and a confidence threshold of 0.1.

840    To demonstrate the ability of the TranSPHIRE feedback loop to selectively pick

841    particles of one conformational state, Cinderella was trained on pre-existing 2D class

842    averages of the pore state as instances of "good" classes (318) and 2D class averages

843    of the pre-pore state and contamination as instances of "bad" classes (664). During

844    the feedback rounds crYOLO was trained on the default value of 50 random

845    micrographs that contained particles contributing to classes labeled "good" by

846    Cinderella. 3D reconstructions were computed as described above without applying a

847    mask or symmetry.

848        The progression of the picking performance throughout the feedback rounds

849   was evaluated on a fixed subset of 500 micrographs as described above (Figure 4).

850

851   **Automatic processing of an actomyosin complex data set.**

852   A subset of 100 micrographs of an actomyosin complex with a bound small molecule

853   ligand (publication in preparation) was chosen to demonstrate the processing of

854   filamentous samples and within TranSPHIRE and its suitability for high-throughput

855   ligand screenings. The data set was collected at a Cs-corrected Titan Krios (FEI

856   Thermo Fisher) microscope equipped with an X-FEG and operated at 300kV using

857   EPU (FEI Thermo Fisher). Equally dosed frames with a pixel size of 0.56 Å/pixel were

858   collected with a K2 Summit (super resolution mode, Gatan) direct electron detector

859   in combination with a GIF quantum-energy filter set to a filter width of 20 eV.  Each

860   movie contains 40 frames and a total electron dose of 81.2 e/Å$^2$.

861   Processing in TranSPHIRE was performed as described above with five internal

862   feedback rounds to optimize the crYOLO picking model. Within the pipeline, movies

863   were drift corrected, dose weighted and binned to a pixel size of 1.10 Å/px by

864   MotionCor2 [14] deactivating patch alignment and CTFFIND4 [16] fitted the CTF between

865   5 Å and 30 Å with an Cs value of 0.001.

866        As the crYOLO general model does not include filamentous data it cannot be

867   readily applied to this data set. Instead a new crYOLO general model specific for actin

868   filaments was trained. The training data consisted of multiple actin data sets

869   collected within our group, but did not include any data of an actomyosin complex or

870   other actin complexes. Considering the significant optical difference of actin and

871   actomyosin filaments (also see Figure 5), picking with the general actin crYOLO model

872   mimics the processing of a previously unseen filamentous protein.

36

873 During the feedback rounds the crYOLO picking threshold was set to 0.1 and

874 the anchor size to the estimated box size of 320 pixels. Furthermore, the filament

875 width was set to 100 px and the box distance to 25 px (equivalent to one helical rise

876 of 27.5 Å). Only filaments consisting of at least six segments were considered. After

877 the final feedback round, the picking threshold value was adjusted to the crYOLO

878 default value of 0.3, as the threshold optimization procedure of crYOLO does not

879 support filaments. After each particle picking step, particles were automatically

880 extracted using SPHIRE sp_window.py with a box size of 320 pixels and a filament

881 width of 100 pixels. The subsequent 2D classification was performed using a GPU

882 accelerated version of the SPHIRE ISAC2 algorithm asking for 30-50 particles per

883 class. The feedback loop was run with the default particle batch size of 20,000 (for

884 details see above and Supplementary Figure 3).

885 The produced 2D class averages were subjected to an automatic 2D class

886 selection using our deep learning tool Cinderella and a confidence threshold of 0.1.

887 As filamentous data differ strongly from the data used to train the general model of

888 Cinderella, a new model was trained based on the 2D class averages produced in the

889 initial feedback round combined with previously selected class averages of actin only

890 data sets. During the feedback rounds crYOLO was trained on the default value of 50

891 random micrographs that contained particles contributing to classes labeled "good"

892 by Cinderella.

893 An initial 3D reference was created from a deposited actomyosin atomic

894 model (PDB:5JLH) [38]. The 3D refinement was performed using SPHIRE

895 sp_meridien_alpha.py, an open alpha version of helical processing in SPHIRE, with a

896 particle radius of 144 px (~45% of the box size), a filament width of 100 px and a

897 helical rise of 27.5 Å. While projection parameters are restrained according to the

898 helical parameters e.g. the shift along the filament axis is restricted to half of the rise,

37

899    no helical symmetry is applied and therefore does not need to be determined

900    beforehand. To avoid artifacts due to the contact of the filament to the edges of the

901    box, a soft 3D mask covering 85% percent of the filament was applied during the

902    refinement.

903    The progression of the picking performance throughout the feedback rounds was

904    evaluated on a fixed subset of 100 micrographs as described above (Figure 5).

905

## Acknowledgements

912

## Author contributions

914    **Conceptualization**: M.S., T.W., C.G. and S.R.;

915    **Software - TranSPHIRE:** M.S.;

916    **Software - GPU ISAC:** F.S.;

917    **Software - Cinderella:** T.W.;

918    **Formal Analysis:** M.S., T.W., S.P.;

919    **Writing – Original Draft:** F.S.;

920    **Writing – Review & Editing:** F.S., M.S., S.P., T.W., C.G., S.R.;

921    **Funding Acquisition:** S.R.

922

923

924 **Competing interests**

925 The authors declare no competing interests.

926

927 **Data availability**

928 The movies processed in this manuscript are subsets of data sets that are published

929 (Tc holotoxin [29]) or will be published elsewhere (TRPC4 and actomyosin) and are

930 available from the corresponding author upon reasonable request.

931

932 **Code availability**

933 TranSPHIRE is open-source and can be downloaded free of charge

934 (https://github.com/MPI-Dortmund/transphire).

935

936 **References**

937 1.    Nogales, E. The development of cryo-EM into a mainstream  structural biology

938       technique. *Nat. Methods* **13,** 24–27 (2016).

939 2.    Method of the Year 2015. *Nature Publishing Group* 1–1 (2015).

940       doi:10.1038/nmeth.3730

941 3.    Merino, F. & Raunser, S. Cryo-EM as a tool for structure-based drug

942       development. *Angewandte Chemie* (2016). doi:10.1002/ange.201608432

943 4.    Schenk, A. D., Cavadini, S., Thomä, N. H. & Genoud, C. Live Analysis and

944       Reconstruction of Single-Particle Cryo-Electron Microscopy Data with

945       CryoFLARE. *J Chem Inf Model* acs.jcim.9b01102 (2020).

946       doi:10.1021/acs.jcim.9b01102

947 5.    Zivanov, J. *et al.* New tools for automated high-resolution cryo-EM structure

948       determination in RELION-3. *Elife* **7,** (2018).

39

949   6.    Tegunov, D. & Cramer, P. Real-time cryo-electron microscopy data

950         preprocessing with Warp. *Nat. Methods* **16,** 1146–1152 (2019).

951   7.    Maluenda, D. *et al.* Flexible workflows for on-the-fly electron-microscopy

952         single-particle image processing using Scipion. *Acta Crystallogr D Struct Biol*

953         **75,** 882–894 (2019).

954   8.    Biyani, N. *et al.* Focus: The interface between data collection and data

955         processing in cryo-EM. *J. Struct. Biol.* **198,** 124–133 (2017).

956   9.    Punjani, A., Rubinstein, J. L., Fleet, D. J. & Brubaker, M. A. cryoSPARC:

957         algorithms for rapid unsupervised cryo-EM structure determination. *Nat.*

958         *Methods* (2017). doi:10.1038/nmeth.4169

959   10.   Lander, G. C. *et al.* Appion: an integrated, database-driven pipeline to

960         facilitate EM image processing. *J. Struct. Biol.* **166,** 95–102 (2009).

961   11.   Suloway, C. *et al.* Automated molecular microscopy: the new Leginon system.

962         *J. Struct. Biol.* **151,** 41–60 (2005).

963   12.   Wagner, T. *et al.* SPHIRE-crYOLO is a fast and accurate fully automated particle

964         picker for cryo-EM. *Commun Biol* **2,** 218 (2019).

965   13.   Moriya, T. *et al.* High-resolution Single Particle Analysis from Electron Cryo-

966         microscopy Images Using SPHIRE. *J Vis Exp* e55448–e55448 (2017).

967         doi:10.3791/55448

968   14.   Zheng, S. Q. *et al.* MotionCor2: anisotropic correction of beam-induced

969         motion for improved cryo-electron microscopy. *Nat. Methods* (2017).

970         doi:10.1038/nmeth.4193

971   15.   Grant, T. & Grigorieff, N. Measuring the optimal exposure for single particle

972         cryo-EM using a 2.6 Å reconstruction of rotavirus VP6. *Elife* **4,** e06980 (2015).

973   16.   Rohou, A. & Grigorieff, N. CTFFIND4: Fast and accurate defocus estimation

974         from electron micrographs. *J. Struct. Biol.* **192,** 216–221 (2015).

975   17.   Penczek, P. A. *et al.* CTER-rapid estimation of CTF parameters with error

976          assessment. *Ultramicroscopy* **140,** 9–19 (2014).

977   18.   Zhang, K. Gctf: Real-time CTF determination and correction. *J. Struct. Biol.*

978          **193,** 1–12 (2016).

979   19.   Yang, Z., Fang, J., Chittuluru, J., Asturias, F. J. & Penczek, P. A. Iterative stable

980          alignment and clustering of 2D transmission electron microscope images.

981          *Structure* **20,** 237–247 (2012).

982   20.   Wagner, T. Cinderella. (2019). doi:10.5281/zenodo.3672421

983   21.   Hohn, M. *et al.* SPARX, a new environment for Cryo-EM image processing. *J.*

984          *Struct. Biol.* **157,** 47–55 (2007).

985   22.   Bepler, T. *et al.* Positive-unlabeled convolutional neural networks for particle

986          picking in cryo-electron micrographs. *Nat. Methods* **16,** 1153–1160 (2019).

987   23.   Wang, F. *et al.* DeepPicker: A deep learning approach for fully automated

988          particle picking in cryo-EM. *J. Struct. Biol.* **195,** 325–336 (2016).

989   24.   Behrmann, E. *et al.* Real-space processing of helical filaments in SPARX. *J.*

990          *Struct. Biol.* **177,** 302–313 (2012).

991   25.   He, S. & Scheres, S. Helical reconstruction in RELION. 1–27 (2016).

992          doi:10.1101/095034

993   26.   Rohou, A. & Grigorieff, N. Frealix: model-based refinement of helical filament

994          structures from electron micrographs. *J. Struct. Biol.* **186,** 234–244 (2014).

995   27.   Egelman, E. H. The iterative helical real space reconstruction method:

996          surmounting the problems posed by real polymers. *J. Struct. Biol.* **157,** 83–94

997          (2007).

998   28.   Wagner, T. *et al.* Two particle picking procedures for filamentous proteins:

999          SPHIRE-crYOLO filament mode and SPHIRE-STRIPER. **2,** 218–23 (2020).

29. Roderer, D., Hofnagel, O., Benz, R. & Raunser, S. Structure of a Tc holotoxin pore provides insights into the translocation mechanism. *Proc. Natl. Acad. Sci. U.S.A.* **116,** 23083–23090 (2019).

30. Kremer, J. R., Mastronarde, D. N. & McIntosh, J. R. Computer visualization of three-dimensional image data using IMOD. *J. Struct. Biol.* **116,** 71–76 (1996).

31. Sokolova, M. & Lapalme, G. A systematic analysis of performance measures for classification tasks. *Information Processing and Management* **45,** 427–437 (2009).

32. Iudin, A., Korir, P. K., Salavert-Torres, J., Kleywegt, G. J. & Patwardhan, A. EMPIAR: a public archive for raw electron microscopy image data. *Nat. Methods* **13,** 387–388 (2016).

33. Mirjalili, S., Song Dong, J., Sadiq, A. S. & Faris, H. in *Nature-Inspired Optimizers: Theories, Literature Reviews and Applications* (eds. Mirjalili, S., Song Dong, J. & Lewis, A.) 69–85 (Springer International Publishing, 2020).

34. Elmlund, H., Elmlund, D. & Bengio, S. PRIME: probabilistic initial 3D model generation for single-particle cryo-electron microscopy. *Structure* **21,** 1299–1306 (2013).

35. Henderson, R. *et al.* Outcome of the first electron microscopy validation task force meeting. in **20,** 205–214 (2012).

36. Vinayagam, D. *et al.* Electron cryo-microscopy structure of the canonical TRPC4 ion channel. *Elife* **7,** 213 (2018).

37. Zivanov, J., Nakane, T. & Scheres, S. H. W. Estimation of high-order aberrations and anisotropic magnification from cryo-EM data sets in RELION-3.1. *IUCrJ* **7,** 253–267 (2020).

1024    38.    Ecken, von der, J., Heissler, S. M., Pathan-Chhatbar, S., Manstein, D. J. &

1025         Raunser, S. Cryo-EM structure of a human cytoplasmic actomyosin complex at

1026         near-atomic resolution. *Nature* **534,** 724–728 (2016).

1027