ACCEPTED MANUSCRIPT

eLIFE

Low cost, high performance processing of single particle cryo-electron microscopy data in the cloud

Michael A Cianfrocco, Andres E Leschziner

This PDF is the version of the article that was accepted for publication after peer review. Fully formatted HTML, PDF, and XML versions will be made available after technical processing, editing, and proofing.

Stay current on the latest in life science and biomedical research from eLife.
Sign up for alerts at elife.elifesciences.org

1

2

3

4

5

6

7

8

9

10   **Low cost, high performance processing of single particle cryo-electron microscopy**

11   **data in the cloud**

12   Michael A. Cianfrocco[1,2*] and Andres E. Leschziner[1]

13

14   [1]Department of Molecular and Cellular Biology, Harvard University, Cambridge, MA, 02138 USA

15   [2]Department of Cell Biology, Harvard Medical School, Boston, MA, 02115, USA

16   *Corresponding author: mcianfrocco@fas.harvard.edu

17

18

19

20

21

22

23

24

## ABSTRACT

The advent of a new generation of electron microscopes and direct electron detectors has realized the potential of single particle cryo-electron microscopy (cryo-EM) as a technique to generate high-resolution structures. Calculating these structures requires high performance computing clusters, a resource that may be limiting to many likely cryo-EM users. To address this limitation and facilitate the spread of cryo-EM, we developed a publicly available 'off-the-shelf' computing environment on Amazon's elastic cloud computing infrastructure. This environment provides users with single particle cryo-EM software packages and the ability to create computing clusters with 16 to 480+ CPUs. We tested our computing environment using a publicly available 80S yeast ribosome dataset and estimate that laboratories could determine high-resolution cryo-EM structures for $50 to $1,500 per structure within a timeframe comparable to local clusters. Our analysis shows that Amazon's cloud computing environment may offer a viable computing environment for cryo-EM.

**Introduction**

Cryo-electron microscopy (cryo-EM) has long served as an important tool to provide structural insights into biological samples. Recent advances in cryo-EM data collection and analysis, however, have transformed single particle cryo-EM (Bai et al., 2015; Kuhlbrandt, 2014), allowing it to achieve resolutions better than 5 Å for samples ranging in molecular weight from the 4 MDa eukaryotic ribosome (Bai et al., 2013) to the 170 kDa membrane protein γ-secretase (Lu et al., 2014). These high-resolution structures are the result of a new generation of cameras that detect electrons directly without the need of a scintillator, which results in a dramatic increase in the signal-to-noise ratio relative to CCD cameras, the previous most commonly used device (McMullan et al., 2009). In addition to direct electron detection, the high frame rate of these cameras allows each image to be recorded as a 'movie', dividing it into multiple frames. These fractionated images can be used to correct for sample movement during the exposure, further increasing the quality of the cryo-EM images (Campbell et al., 2012; Li et al., 2013; Scheres, 2014).

In addition to these technological developments in the detectors, improvements in computer software packages have played an equally important role in moving cryo-EM into the high-resolution era. Atomic or near-atomic structures have been obtained with software packages such as EMAN2 (Tang et al., 2007), Sparx (Hohn et al., 2007), FREALIGN (Grigorieff, 2007), Spider (Frank et al., 1996), and Relion (Scheres, 2012, 2014). In general, obtaining these structures involved computational approaches that sorted out the data into homogenous classes that could then be refined to high resolution.

65    While these advances in microscopy and analysis have been essential for the

66    recent breakthroughs in cryo-EM, their implementation is computationally intensive and

67    requires high-performance computing clusters. A recent survey of high-resolution single

68    particle cryo-EM structures showed that refinement of these structures required

69    processing times in excess of 1,000 CPU-hours (Scheres, 2014). Therefore,

70    computational time (i.e. access to high-performance clusters) may represent a bottleneck

71    to determining high-resolution structures by single particle cryo-EM.

72    In order to address this limitation, we explored the possibility of using Amazon's

73    elastic cloud computing (EC2) for processing cryo-EM data. To help others take

74    advantage of this resource, we have created a publically available 'off-the-shelf' software

75    environment that allows new users to start up a cluster of Amazon CPUs preinstalled with

76    cryo-EM software and we have used it to test the performance of Amazon's EC2

77    platform. We were able to determine a 4.6 Å structure of the 80S ribosome using a

78    published dataset (Bai et al., 2013) for an overall cost of $100 USD within a timeframe

79    comparable to that of a local cluster. Given the range of prices for accessing Amazon

80    CPUs (users can bid for significantly reduced costs) and the accessibility statistics, we

81    estimate that typical cryo-EM structures can be determined for $50 - $1,500 per structure.

82

83    **Elastic cloud computing through Amazon Web Services**

84    Amazon Web Services (AWS) is a division of Amazon that offers a variety of

85    cloud-based solutions for website hosting and high-performance computing, amongst

86    other services. Many different types of privately held companies take advantage of

87    Amazon's computing infrastructure because of its affordability, flexibility, and security.

88    Of note, global biotechnology companies such as Novartis (AWS, 2014a), Bristol-Myers-

89    Squibb (AWS, 2013), and Pfizer (AWS, 2014b) have utilized the computing power of

90    Amazon for scientific data processing. Many academic researchers have also begun to

91    use Amazon's EC2 resources for analyzing datasets from super-resolution light

92    microscopy (Hu et al., 2013), genomics (Krampis et al., 2012; Yazar et al., 2014), and

93    proteomics (Mohammed et al., 2012; Trudgian and Mirzaei, 2012).

94         The overall workflow starts with users logging into a virtual machine ('instance')

95    on AWS (Figure 1). AWS offers a variety of instance types that have been configured for

96    different computing tasks. For example, instances have been optimized for computing

97    performance, GPU-based calculations, or memory-intensive calculations. After logging

98    onto an instance, storage drives are mounted onto it, allowing data, which can be

99    encrypted for security, to be transferred onto the storage drives (Figure 1).

100        While users can utilize a single instance for calculations, the maximum number of

101   CPU cores per instance is 18. Therefore, creating a computing cluster with a larger

102   number of CPUs on AWS requires additional steps. The Software Tools for Academics

103   and Researchers (STAR) group at the Massachusetts Institute of Technology developed a

104   straightforward package that allows users to group individual AWS instances into a

105   cluster. The STARcluster program is a python-based, open source package that

106   automatically creates a cluster preconfigured with the necessary software to manage a

107   computer cluster (Ivica et al., 2009). This package allows users to specify the number of

108   instances to be included in the clusters as well as the instance type. By taking advantage

109   of this tool, private clusters can be built with sizes ranging from 16 to 480 CPUs (Figure

110   1).

111

**Global availability of spot instances on Amazon EC2**

While Amazon provides dedicated access to instances through 'on-demand'

reservations, there are 'spot instances' that are 80-90% cheaper than the on-demand price.

Spot instances are unused instances within Amazon EC2 that are open for competitive

bidding, where users gain access to them by making offers above the current minimum

bid. This means that while the on-demand rate for high-memory, 16-CPU instances

(called "r3.8xlarge") is $2.80/hr, spot instance prices can be as low as $0.25-$0.35/hr.

In order to determine if spot instances offer a consistent reduction in price, we

analyzed the global availability of r3.8xlarge spot instances. Currently, Amazon has 9

regions worldwide within 7 countries: US-East-1 (United States), US-West-1 (United

States), US-West-2 (United States), SA-East-1 (Brazil), EU-Central-1 (Germany), EU-

West-1 (Ireland), AP-Northeast-1 (Japan), AP-Southeast-1 (Singapore), and AP-

Southeast-2 (Australia). For each region, we retrieved spot instance prices for r3.8xlarge

instances over the past 3 months and analyzed the time they spent at prices below $0.35 -

$0.65/hr (corresponding to discounts of 87.5% to 76.8% over the full on-demand rate of

$2.80/hr) (Figure 2 and Figure 2-figure supplement 1). This analysis revealed that,

globally, 49.8% of r3.8xlarge instances were below $0.35/hr, 12.5% the on-demand price

(Figure 2). For $0.65/hr, 76.5% below full price, one could access 82.2% of the global

r3.8xlarge spot instances. These data indicate that spot instances provide dependable,

cost-effective access to Amazon's computing resources.

132

**Performance analysis of Amazon's EC2 environment with a 80S yeast ribosome**

**dataset**

To test the performance of Amazon's EC2 environment, we analyzed a previously

published 80S *Saccharomyces cerevisiae* ribosome dataset (Bai et al., 2013) (EMPIAR

10002) on a 128 CPU cluster (8 x 16 CPUs; using the r3.8xlarge instance). After

extracting 62,022 particles, we performed 2D classification within Relion. Subsequent

3D classification of the particles into 4 classes revealed that two classes adopted a similar

structural state. We merged those two classes and used the associated particles to carry

out a 3D refinement in Relion—we were able to obtain a structure with an overall

resolution of 4.6 Å (Figure 3A - C).

This structure, whose generation included particle picking, CTF estimation, 2D

and 3D classification, and refinement, cost us $99.64 on Amazon's EC2 environment.

This cost was achieved by bidding on spot instances for particle picking (m1.small at

$0.02/hr), 2D classification (STARcluster of r3.8xlarge instances at $0.65/hr), and 3D

classification and refinement (STARcluster of r3.8xlarge instances at $0.65/hr). Thus,

even though obtaining this structure required 1,266 total CPU-hours, Amazon's EC2

computing infrastructure provided the necessary resources to calculate it to near-atomic

resolution at a reasonable price.

To further test the performance of Amazon instances, we carried out 3D

classification and refinement on a variety of STARcluster configurations using Relion.

As before, we ran our tests on clusters of r3.8xlarge high-memory instances (256 GiB

RAM and 16 CPUs per instance). Comparing performance across cluster sizes showed

that 256 CPUs had the fastest overall time and the highest speedup relative to a single

156 CPU for both 3D classification and refinement (Figure 4A,B). However, cluster sizes of

157 128 and 64 CPUs were the most cost effective for 3D classification and refinement,

158 respectively, as these were the cluster configurations where the speedup per dollar

159 reached a maximum (Figure 4C). Importantly, the average time required to boot up these

160 STARclusters was ≤ 10 minutes for all cluster sizes (Figure 4D) and, once booted up, the

161 clusters do not have any associated job wait times. Therefore, these tests showed that

162 Amazon's EC2 infrastructure was amenable to the analysis of single particle cryo-EM

163 data using Relion over a range of STARcluster sizes.

164 　　　　From our analysis of the 80S yeast ribosome, we extrapolated the processing

165 times and combined them with previously published 3D refinement times to estimate

166 typical costs on Amazon's EC2. First, we estimated the cost for 3D refinement in Relion

167 for previously published structures (Supplementary File 2A)—these calculated costs

168 ranged from $12.65 to $379.03 per structure, depending on the spot instance price and

169 required CPU-hours. We then combined these data with conservative estimates for

170 particle picking, CTF estimation, particle extraction, 2D and 3D classification to predict

171 the overall cost of structure determination on Amazon's EC2 (Supplementary File 2B).

172 From these considerations, we estimated that published structures could be determined

173 using Amazon's EC2 environment at costs of $50 - $1,500 per structure (Supplementary

174 File 2B).

175

176 **EM-packages-in-the-Cloud: A pre-configured software environment for single-**

177 **particle cryo-EM image analysis**

178        Given the success we had in analyzing cryo-EM data on Amazon's EC2 at an

179    affordable price and within a reasonable timeframe, we have made our software

180    environment publicly available as an 'Amazon Machine Image' (AMI), under the name

181    'EM-packages-in-the-Cloud-v3.93.' The EM-packages-in-the-Cloud-v3.93 AMI provides

182    the software environment necessary for analyzing data on a single instance, and is

183    preconfigured with STARcluster software.  The EM-packages-in-the-Cloud-v3.93 AMI

184    has the following cryo-EM software packages installed: Relion (Scheres, 2012, 2014),

185    FREALIGN (Grigorieff, 2007), EMAN2 (Tang et al., 2007), Sparx (Hohn et al., 2007),

186    Spider (Frank et al., 1996), EMAN (Ludtke et al., 1999), and XMIPP (Sorzano et al.,

187    2004). In addition to this AMI that is capable of running on a single instance, we have

188    also made available a second AMI – EM-packages-in-the-Cloud-Node-v3.1 – that

189    provides users with the same software packages as described above, but can set up and

190    run within a cluster of multiple EC2 instances. These two publicly available AMIs allows

191    users to boot up a cluster to analyze cryo-EM data in a few short steps. The protocols

192    describing this can be found as a PDF (Supplementary File 1) or on a Google site that is

193    being launched in conjunction with this article: http://goo.gl/AIwZJz. In addition to

194    detailed instructions, the site includes a help forum to facilitate a conversation on cloud

195    computing for single particle cryo-EM.

196

197
198    **Cloud computing as a tool to facilitate high-resolution cryo-EM**
199
200        Recent advances in single particle cryo-EM have drawn the interest of the broader

201    scientific community. In addition to technical advances in electron optics, the new direct

202    electron detectors and data analysis software have dramatically improved the resolutions

203    that can be achieved for a variety of structural targets. In contrast to the other high-

204    resolution techniques (X-ray crystallography, NMR), structure determination by cryo-EM

205    is extremely computationally intensive. The publicly available 'EM-packages-in-the-

206    Cloud' environment we have presented and characterized here will help remove some of

207    the limitations imposed by these computational requirements.

208         We believe that cloud-based approaches have the potential to impact the future of

209    cryo-EM image processing in two fronts: 1) new cryo-EM users or laboratories will have

210    immediate access to a high performance cluster, and 2) existing labs may use this

211    resource to increase their productivity. As the number of laboratories using cryo-EM

212    increases, and as existing laboratories begin to pursue high-resolution cryo-EM, gaining

213    immediate access to a high performance cluster may become difficult. For instance, while

214    there are government-funded high performance clusters in the United States (e.g. XSEDE

215    STAMPEDE), it may take up to a month for a user application to be reviewed  (Rogelio

216    Hernandez-Lopez, personal communication). Assuming that the application is approved,

217    these clusters may not have appropriate software installed, which further delays data

218    processing. Finally, the user will have a set limit for the number of CPU hours available

219    per project, requiring a new application to be submitted to access the cluster again. All of

220    these problems can be circumvented by using Amazon's EC2 infrastructure, which

221    provides immediate, cost-effective access to hundreds of CPUs with no geographic

222    restrictions.

223         The power of cloud-based solutions to alleviate the computational burden

224    associated with cryo-EM data processing stems from its high-degree of scalability and

225    reasonable cost. By minimizing computational time and increasing global accessibility,

226 high-performance cloud computing may help usher in the era when high-resolution cryo-

227 EM becomes a routine structural biology tool.

228
229 **Materials and Methods**
230

231 *Global availability of spot instances*

232 Global spot instance prices were retrieved from the 90-day period from January 1, 2015

233 to April 1, 2015 using the Amazon Command Line Tools command *ec2-describe-spot-*

234 *price-history*. Retrieval of spot instance prices for all regions was implemented

235 automatically in a custom python program *get_spot_histories_all_regions_all_zones.py*.

236 From these spot instance prices, the percentage time spent below given prices was

237 calculated using *measure_time_at_spotPrice.py*, where the cumulative time of spot

238 instances below a given price were divided by the total time (90 days). Both programs

239 can be found in the Github repository mcianfrocco/Cianfrocco-and-Leschziner-

240 EMCloudProcessing.

241

242 *Setting up a cluster on Amazon EC2 with spot instances*

243 In order to minimize costs, STARclusters were assembled from 'spot instances,' which

244 are unused instances that can be reserved through a bidding process. The spot instances

245 are different from 'on-demand' instances: on-demand instances provide users with

246 guaranteed access while spot instances are reserved until there is a higher bid, at which

247 point the user is logged out of the spot instance. When this happens, the MPI-threaded

248 Relion calculation will abort, requiring the user to resubmit the job to the STARcluster

249 and start Relion from the previous iteration. Even if the user is logged out of all instances

250 within a STARcluster, the data is automatically saved within the EBS-backed volumes on

251 Amazon EC2.

252

253 *CPUs vs. vCPUs*

254 In selecting an instance type, new users should be aware of the differences between CPUs

255 and vCPUs on Amazon's EC2 network. Namely, that there are two vCPUs per physical

256 CPU on Amazon. This means that while r3.8xlarge instances have 32 vCPUs, there are

257 actually only 16 physical CPU cores in each instance, with each CPU having two

258 hyperthreads. Practically, this means that Amazon's instances have higher performance

259 than a 16 CPU machine and less performance than a 32 CPU machine. To account for

260 this difference, all numbers reported here were CPU numbers that were converted from

261 vCPUs: 1 CPU = 2 vCPUs.

262

263 *Image processing*

264 Micrographs from the 80S *Saccharomyces cerevisiae* ribosome dataset (Bai et al., 2013)

265 were downloaded from the EMPIAR database for electron microscopy data (EMPIAR

266 10002). The SWARM feature of EMAN2 (Tang et al., 2007) was used to pick particles

267 semi-automatically. Micrograph defocus was estimated using CTFFIND3 (Mindell and

268 Grigorieff, 2003). The resulting particle coordinates and defocus information were used

269 for particle extraction by Relion-v1.3 (Scheres, 2012, 2014). The particle stacks and

270 associated data files were then uploaded to an elastic block storage volume on Amazon's

271 EC2 processing environment at a speed of 10 MB/sec (24 minutes total upload time).

272    After 2D classification in Relion, 3D classification was performed on 62,022 80S

273    Ribosome particles (1.77 Å/pixel), also in Relion. These were classified into 4 groups

274    (T=4) for 13 iterations using a ribosome map downloaded from the Electron Microscopy

275    Data Bank (EMDB-1780) that was low pass filtered to 60 Å. Further 3D classification

276    using a local search of 10° and an angular sampling of 1.8° continued for 13 iterations. At

277    this point, two classes were identified as belonging to the same structural state and were

278    selected for high-resolution refinement (32,533 particles). Refinement of these selected

279    particles continued for 31 iterations using *3D auto-refine* in Relion. The final resolution

280    was determined to be 4.6 Å using *Post process* in Relion, applying a mask to the merged

281    half volumes and a negative B-factor of -116 Å$^2$.

282

283    *Performance Analysis*

284    80S ribosome data were reanalyzed on clusters of increasing size using both 3D

285    classification and 3D refinement. The time points collected involved running 3D

286    classification for 2 rounds and 3D refinement for 6 rounds, using the same number of

287    particles and box sizes listed above: 62,022 particles for classification and 32,533

288    particles for refinement with box sizes of 240 x 240 pixels. The Relion commands were

289    identical to the commands used above and the calculations were terminated after the

290    specified iteration.

291    From these time points, the speedup of each cluster size was calculated relative to

292    a single CPU. Speedup (*S*) was calculated as:

$$S = \frac{\text{Calculation time for 1 CPU}}{\text{Calculation time for } x \text{ CPUs}}$$

293

294    The measured speedup values were then compared to the speedup expected for a

295    perfectly parallel algorithm ($P = 1$) using Amdahl's law (Amdahl, 1967):

296
$$S = \frac{1}{(1-P) + \frac{1}{n}(P)} = \frac{1}{(1-1) + \frac{1}{n}(1)} = n$$

297    Where $P$ is the fraction of an algorithm that is parallel and $n$ is the number of processors.

298    The calculation times for 3D classification on a single CPU were obtained by using 1

299    CPU on a 16 CPU r3.8xlarge instance. For calculating a 3D refinement on a single CPU,

300    (or two vCPUs), the refinement was run on 4 vCPUs and then converted to a single CPU

301    (or two vCPUs) by multiplying the calculation time by 2. For cost analysis, the measured

302    speedup was divided by the cost to run the job on spot instances of r3.8xlarge at a price

303    of $0.35/hr. Cluster boot up times were calculated from the elapsed time between

304    submitting the STARcluster command and the STARcluster fully booting up.

305

306    **Data Accession Information**

307    Further information regarding 'EM-Packages-in-the-Cloud' can be found in

308    Supplementary File 1 and at an associated Google Site: http://goo.gl/AIwZJz. The final

309    80S yeast ribosome structure at 4.6 Å has been submitted to the EM Databank as EMDB

310    2858. A detailed description of global spot instance price analyses and image processing

311    is available at https://github.com/mcianfrocco/Cianfrocco-and-Leschziner-

312    EMCloudProcessing/wiki. Associated computing scripts and data files have been

313    uploaded to Github (https://github.com/mcianfrocco/Cianfrocco-and-Leschziner-

314    EMCloudProcessing) and Dryad Digital Repository

315    (http://datadryad.org/review?doi=doi:10.5061/dryad.9mb54) (Cianfrocco and

316    Leschziner), respectively.

**Competing financial interests statement**

The authors do not have any competing financial interests.

**References**

Amdahl, G.M. (1967). Validity of the single processor approach to achieving large scale computing capabilities. In Proceedings of the April 18-20, 1967, spring joint computer conference (Atlantic City, New Jersey: ACM), pp. 483-485.

AWS: Bristol-Myers Squibb on AWS. Date Accessed: April 20, 2015. (http://aws.amazon.com/solutions/case-studies/bristol-myers-squibb/)

AWS: AWS Case Study: Novartis. Date Accessed: April 20, 2015. (http://aws.amazon.com/solutions/case-studies/novartis/)

AWS: AWS Case Study: Pfizer. Date Accessed: April 20, 2015. (http://aws.amazon.com/solutions/case-studies/pfizer/)

Bai, X.C., Fernandez, I.S., McMullan, G., and Scheres, S.H. (2013). Ribosome structures to near-atomic resolution from thirty thousand cryo-EM particles. eLife *2*, e00461.

Bai, X.C., McMullan, G., and Scheres, S.H. (2015). How cryo-EM is revolutionizing structural biology. Trends in biochemical sciences *40*, 49-57.

Campbell, M.G., Cheng, A., Brilot, A.F., Moeller, A., Lyumkis, D., Veesler, D., Pan, J., Harrison, S.C., Potter, C.S., Carragher, B.*, et al.* (2012). Movies of ice-embedded particles enhance resolution in electron cryo-microscopy. Structure (London, England : 1993) *20*, 1823-1828.

Cianfrocco, M.A., and Leschziner, A.E. Data from: Single particle cryo-electron microscopy image processing in the cloud: High performance at low cost (Dryad Data Repository).

Frank, J., Radermacher, M., Penczek, P., Zhu, J., Li, Y., Ladjadj, M., and Leith, A. (1996). SPIDER and WEB: processing and visualization of images in 3D electron microscopy and related fields. Journal of structural biology *116*, 190-199.

Grigorieff, N. (2007). FREALIGN: high-resolution refinement of single particle structures. Journal of structural biology *157*, 117-125.

Hohn, M., Tang, G., Goodyear, G., Baldwin, P.R., Huang, Z., Penczek, P.A., Yang, C., Glaeser, R.M., Adams, P.D., and Ludtke, S.J. (2007). SPARX, a new environment for Cryo-EM image processing. Journal of structural biology *157*, 47-55.

363 Hu, Y.S., Nan, X., Sengupta, P., Lippincott-Schwartz, J., and Cang, H. (2013).
364 Accelerating 3B single-molecule super-resolution microscopy with cloud computing.
365 Nature methods *10*, 96-97.

366 Ivica, C., Riley, J.T., and Shubert, C. (2009). StarHPC - Teaching parallel programming
367 within elastic compute cloud. Paper presented at: Information Technology Interfaces,
368 2009 ITI '09 Proceedings of the ITI 2009 31st International Conference on.

369 Krampis, K., Booth, T., Chapman, B., Tiwari, B., Bicak, M., Field, D., and Nelson, K.E.
370 (2012). Cloud BioLinux: pre-configured and on-demand bioinformatics computing for
371 the genomics community. BMC bioinformatics *13*, 42.

372 Kuhlbrandt, W. (2014). Cryo-EM enters a new era. eLife *3*, e03678.

373 Li, X., Mooney, P., Zheng, S., Booth, C.R., Braunfeld, M.B., Gubbens, S., Agard, D.A.,
374 and Cheng, Y. (2013). Electron counting and beam-induced motion correction enable
375 near-atomic-resolution single-particle cryo-EM. Nature methods *10*, 584-590.

376 Lu, P., Bai, X.C., Ma, D., Xie, T., Yan, C., Sun, L., Yang, G., Zhao, Y., Zhou, R.,
377 Scheres, S.H.*, et al.* (2014). Three-dimensional structure of human gamma-secretase.
378 Nature *512*, 166-170.

379 Ludtke, S.J., Baldwin, P.R., and Chiu, W. (1999). EMAN: semiautomated software for
380 high-resolution single-particle reconstructions. Journal of structural biology *128*, 82-97.

381 McMullan, G., Faruqi, A.R., Henderson, R., Guerrini, N., Turchetta, R., Jacobs, A., and
382 van Hoften, G. (2009). Experimental observation of the improvement in MTF from
383 backthinning a CMOS direct electron detector. Ultramicroscopy *109*, 1144-1147.

384 Mindell, J.A., and Grigorieff, N. (2003). Accurate determination of local defocus and
385 specimen tilt in electron microscopy. Journal of structural biology *142*, 334-347.

386 Mohammed, Y., Mostovenko, E., Henneman, A.A., Marissen, R.J., Deelder, A.M., and
387 Palmblad, M. (2012). Cloud parallel processing of tandem mass spectrometry based
388 proteomics data. Journal of proteome research *11*, 5101-5108.

389 Pettersen, E.F., Goddard, T.D., Huang, C.C., Couch, G.S., Greenblatt, D.M., Meng, E.C.,
390 and Ferrin, T.E. (2004). UCSF Chimera--a visualization system for exploratory research
391 and analysis. Journal of computational chemistry *25*, 1605-1612.

392 Scheres, S.H. (2012). RELION: implementation of a Bayesian approach to cryo-EM
393 structure determination. Journal of structural biology *180*, 519-530.

394    Scheres, S.H. (2014). Beam-induced motion correction for sub-megadalton cryo-EM
395    particles. eLife *3*, e03665.


396    Sorzano, C.O., Marabini, R., Velazquez-Muriel, J., Bilbao-Castro, J.R., Scheres, S.H.,
397    Carazo, J.M., and Pascual-Montano, A. (2004). XMIPP: a new generation of an open-
398    source image processing package for electron microscopy. Journal of structural biology
399    *148*, 194-204.


400    Tang, G., Peng, L., Baldwin, P.R., Mann, D.S., Jiang, W., Rees, I., and Ludtke, S.J.
401    (2007). EMAN2: an extensible image processing suite for electron microscopy. Journal
402    of structural biology *157*, 38-46.


403    Trudgian, D.C., and Mirzaei, H. (2012). Cloud CPFP: a shotgun proteomics data analysis
404    pipeline using cloud and high performance computing. Journal of proteome research *11*,
405    6282-6290.


406    Yazar, S., Gooden, G.E., Mackey, D.A., and Hewitt, A.W. (2014). Benchmarking
407    undedicated cloud computing providers for analysis of genomic datasets. PloS one *9*,
408    e108490.
409
410
411
412

**Figure legends and tables**

**Figure 1: Workflow for analyzing cryo-EM data on Amazon's cloud computing infrastructure** After collecting cryo-EM data (Step 1), particles are extracted from the micrographs and prepared for further analysis (Step 2). After logging into an 'instance' (Step 3), data are uploaded to a storage server (elastic block storage) (Step 4). At this point, STARcluster can be configured to launch a cluster of 2 – 30 instances that is mounted with the data from the storage volume (Step 5). A detailed protocol can be found at an accompanying Google site: http://goo.gl/AIwZJz.

**Figure 2: Global availability of Amazon r3.8xlarge spot instances.** Shown is the average percentage time spent by the r3.8xlarge type of instance when the current spot instance price was less than the queried price. The data are averaged over all Amazon's regions worldwide (except for SA-East-1, which does not offer r3.8xlarge instances). Spot instance prices were calculated over a 90-day period from January 1, 2015 – April 1, 2015, where the average is shown +/- the s.e. **Source data:** Figure 2 – Source data 1.

**Figure 2 –figure supplement 1: Availability of virtual machines within regions at specified spot instance prices.** For each Amazon region (excluding SA-East-1, which does not offer r3.8xlarge instances), r3.8xlarge spot instance prices were retrieved for each availability zone, where separate availability zones are shown as separate data points for a given spot instance price. (Note: each region can have different number of availability zones). From the spot instance prices, the percentage time of the spot instances that were spent below the specified spot instance price were calculated. The average value is shown as a solid black line. **Source data:** Figure 2 – Source data 1.

**Figure 3: Cryo-EM structure of 80S ribosome at an overall resolution of 4.6 Å.** (A) Overall view of 80S reconstruction filtered to 4.6 Å while applying a negative B-factor of -116 Å$^2$. (B) Gold standard FSC curve. (C) Selected regions from the 60S subunit. Cryo-EM maps were visualized with UCSF Chimera (Pettersen et al., 2004). **Source data:** Dryad Digital Repository dataset (http://datadryad.org/review?doi=doi:10.5061/dryad.9mb54) (Cianfrocco and Leschziner).

**Figure 4: Relion performance on STARcluster configurations of Amazon instances** (A) Processing times (minutes) for Relion to perform 3D Classification or 3D refinement on 80S ribosome dataset. (B) Speedup for each cluster size relative to a single CPU (black line) shown alongside performance estimate for a perfectly parallel cluster using Amdahl's Law (curve labeled "Theoretical limit"). For cluster sizes ≤ 64 CPUs, Relion exhibits near-perfect performance on STARcluster configurations, while cluster sizes > 64 show that Relion's performance reaches a maximum at 256 CPUs for both 3D classification and 3D refinement. (C) Speedup/Cost is plotted against cluster size, where Speedup/Cost is defined as the speedup observed divided by the cost associated with Amazon's pricing at $0.35/hr/16 CPUs. (D) Average STARcluster boot up time (+/- s.d.) was measured for clusters of increasing size (n = 5). **Source data:** Figure 4 – Source data 1.

**Source data files:**
Figure 2–source data 1: Global spot instance price data from January 1, 2015 to April 1, 2015.
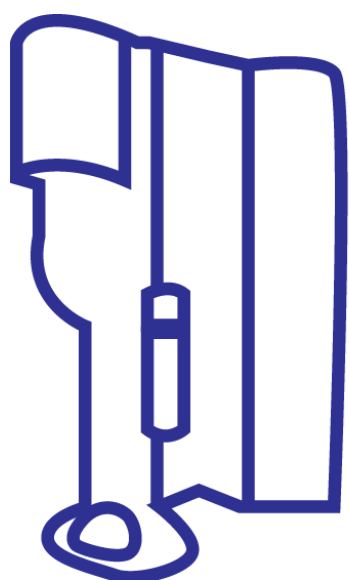
Figure 4–source data 1: Performance analysis statistics for Relion 3D classification and 3D refinement on STARcluster configurations.
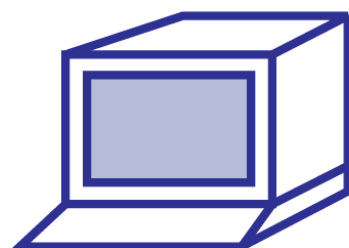
**Supplemental files:**
Supplementary File 1: Step-by-step tutorial describing how to use Amazon's EC2 environment to analyze cryo-EM data.

Supplementary File 2: Comparison of estimated processing times and costs for recent near-atomic cryo-EM structures on Amazon's EC2.

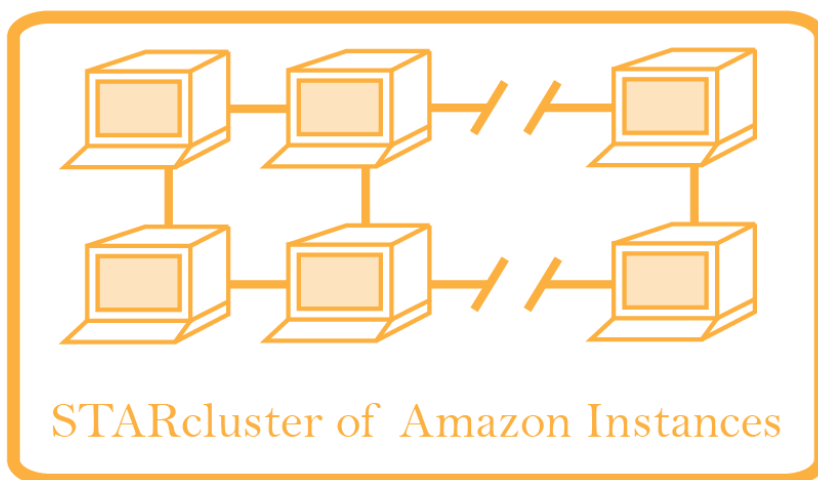Supplementary File 3: Source data for tables in Supplementary File 2

Amazon Elastic Cloud Computing

STARcluster of Amazon Instances

1. Collect cryo-EM data

2. Prepare particles for refinement

3. Log into Amazon instance

4. Upload data onto storage drives

5. Boot up cluster

**A** 60S 40S

**B**

Gold-standard FSC

FSC=0.143

4.6 Å

Resolution (Å)

**C**