

Automated interpretation of Cryo-EM density maps with convolutional neural networks

Philipp Mostosi^{1,2}, Hermann Schindelin¹, Philip Kollmannsberger², Andrea Thorn^{*1}

¹Institute of Structural Biology, Rudolf Virchow Center for Experimental Biomedicine, University of Würzburg, Josef-Schneider-Str. 2, 97080 Würzburg, Germany

²Center for Computational and Theoretical Biology, University of Würzburg, Campus Hubland Nord 32, 97074 Würzburg, Germany

Abstract

Haruspex is a fully convolutional neural network that automatically annotates both protein secondary structure and nucleotides in experimentally derived Cryo-EM maps. The network was trained on a carefully curated dataset of EMDB (Electron Microscopy Data Bank) entries. Haruspex enables users to identify folds and can be used to guide model building as well as validate structures.

In recent years, three-dimensional density maps reconstructed from single particle images obtained by electron cryo-microscopy (Cryo-EM) have reached unprecedented resolution. However, modelling an atomic structure to these maps remains difficult as researchers mostly rely on algorithms developed for crystallographic electron density maps, which are different in both their nature and error distribution. The first step in modelling a reconstruction map - assigning a fold to map regions - can be a major challenge. Parallel to the advances in Cryo-EM in the last decade, deep neural networks achieved remarkable image segmentation capabilities [1], making them the most powerful machine learning approach available. Convolutional neural networks (CNN) combine traditional image analysis with machine learning by cascading layers of trainable convolution filters. CNNs are thus exceptionally well suited for volume annotation and have been successfully applied to biological problems such as breast cancer mitosis recognition [2] and, in conjunction with encoder-decoder architectures, to volumetric data segmentation [3;4].

In this work, we demonstrate that deep neural networks are capable of assigning macromolecular type, i.e. protein or nucleotides, and protein secondary structure elements to experimentally derived Cryo-EM maps. This can be utilized to drastically facilitate model building, to validate existing models and support the placement of known domain folds.

In low-resolution Cryo-EM maps, α -helices can often be discerned as long cylindrical elements. This has been exploited by the program *helixhunter* [5], which searches for prototype helices in reconstruction maps using a cross-correlation strategy. β -Strands are more difficult to identify as they are more variable in shape and therefore require morphological analysis [6]. A combination of these approaches led to the development of *SSEHunter* [7], which uses a density skeleton to detect

37 secondary structures. Deep learning poses an alternative approach and here we demonstrate that
38 resolutions of 4 Å or better, experimental data allow training a well-performing network for a
39 multitude of specimens. Fully convolutional networks [8;3] allow swift segmentation map generation
40 for objects of variable size and we employed a state of the art U-Net-style architecture [3]. The
41 network processed 40^3 voxel segments with a voxel size of 1.0-1.2 Å³ (covering a secondary structure
42 element and its immediate surroundings) to annotate 20^3 voxel cubes (corresponding to the center
43 of the input volume) with four channels containing the probabilities that the voxel is part of an
44 α -helical or β -strand protein secondary structure element, RNA/DNA nucleotide, or other.

45 For network training, we pre-selected EMDB (Electron Microscopy Data Bank [9]) reconstruction
46 maps with resolutions of 4 Å or better. From 576 entries (as of 15/2/2018), we picked 293
47 EMDB/PDB (Protein Data Bank [10]) pairs by three criteria: (1) Good fit between map and model; (2)
48 presence of at least one annotated α -helix or β -sheet; (3) preference of higher resolution maps in
49 case the same authors deposited several instances of the same macromolecular complex. Maps with
50 severe misfits, misalignments, or models without corresponding reconstruction density (and vice
51 versa) were omitted.

52 To generate ground truth data for network training, a python script was implemented to
53 automatically annotate the reconstruction map according to the deposited structural model as
54 α -helical, β -strand, nucleotide or other. The script combined the original annotations from PDBML
55 files with secondary structure identified by DSSP [11] and STRIDE-like extension [12;13] (see
56 Methods). If a secondary structure was identified, all voxels within 3 Å of backbone atoms were
57 annotated accordingly. All voxels with density $\geq 1.0 \sigma$ (standard deviation of the map density
58 distribution) but not within 5 Å of model atoms with density $\geq 1.0 \sigma$ were masked and did not
59 contribute to the training. The voxel size of the reconstruction map was re-scaled to 1.1 Å, if outside
60 [1.0; 1.2] Å. The maps were sliced into a total of 2183 segments á 70^3 Å³ voxels, of which 110
61 segments (5%) were set aside for evaluation during network training. Each segment had to contain
62 at least 100 atoms $\geq 1.0 \sigma$, a backbone mean density of $\geq 3.0 \sigma$, and at least 5% of the total segment
63 volume annotated. The training data were augmented through on-GPU 90° rotations (24
64 possibilities), and by selecting a 40^3 voxel sub-segment at a random position. The network was
65 trained for 40,000 steps with 100 segments employed per step.

66 After training, the network was tested on an independent set of 167 EMDB maps (selected by the
67 same criteria as training data and deposited after February 2018). Virus and ribosome structures
68 were omitted from the test set: viruses' symmetry definition can disagree between map and model
69 and symmetry-averaged maps exhibit particular features; ribosomes are very common and may
70 hence bias the network. For evaluation, we investigated residues with mean backbone densities ≥ 1.0
71 σ and compared the predicted secondary structure on a per-residue basis with the one derived from
72 the deposited PDB model. Using this criterion, the network achieved similar performance on
73 training, evaluation and test data. Over all test maps, there were 74.1% true positives r_p (correctly
74 predicted residues), 18.9% false positives f_p and 4.4% false negatives f_n , resulting in a median recall
75 rate $r_p/(r_p + f_n)^{-1}$ of 94.2% and precision $r_p/(r_p + f_p)^{-1}$ of 79.2%. As a typical example the human
76 ribonuclease P holoenzyme (EMDB entry 9627) illustrates the power of our approach (Fig. 1), which

77 is not only able to accurately predict the RNA vs. protein distribution in this complex but also
78 correctly assigns secondary structure elements in the protein areas with few exceptions. While the
79 high number of false positives was worrying at first, inspection of the test cases revealed that false
80 positives were often elements closely resembling helices, sheets or nucleotides (see Fig. 2). In
81 particular, semi-helical structures, β -hairpin turns and residues belonging to polyproline type II (P_{II})
82 helices [14] were misclassified as α -helical and loosely parallel structures were frequently
83 misclassified as β -strands.

84 Haruspex was trained for resolutions as low as 4 Å, and with the current rate of resolution increase
85 in published maps, by 2021, the average resolution may well be 3.5 Å. Irrespective of this, we will
86 extend our approach to lower resolution data in the future; low resolution experimental maps with a
87 well-matching model, however, are difficult to obtain. This obstacle has previously been faced by Si
88 et al. [15] (*SSELearner*) and Li et al. who developed machine learning approaches for protein
89 secondary structure prediction in Cryo-EM maps (but not nucleotides) [16], and consequently
90 resorted partly to simulated maps generated with pdb2mrc [17]. Simulated maps may lack the error
91 structure and processing artefacts found in experimentally derived reconstruction densities. Si et al.
92 tested their support vector machine on 10 simulated maps of relatively small structures (<40 kDa)
93 and, as available data were still very limited in 2012, only 13 experimental maps paired with
94 individually selected training maps. Haslam et al. [18] used a 3D U-Net, which was trained on 25
95 simulated and 42 experimental maps between 3-9 Å resolution to predict helices and sheets
96 obtaining an $F1$ score $2(\text{recall}^{-1} + \text{precision}^{-1})^{-1}$ between 0.79 and 0.88. However, the network was
97 only tested on six simulated maps and one experimentally derived map. We used a total of 293
98 experimentally derived maps in a semi-automated process to provide a more realistic training
99 environment. The amount of newly released high resolution structures in conjunction with our
100 processing infrastructure allowed us to test our network performance on a representative set of 167
101 unique depositions. In addition, we identify nucleotides, which to our knowledge has not been
102 attempted before. Ribosomes, spliceosomes and polymerases all contain substantial amounts of
103 DNA/RNA nucleotides and are among the most common specimens studied by single-particle Cryo-
104 EM. In addition, β -turns, poly-proline and membrane detergent regions might be desirable additions
105 in future versions of the network.

106 We show that a neural network can be used to automatically distinguish between nucleotides and
107 protein and to assign the two main protein secondary structure elements in experimentally derived
108 Cryo-EM maps. This technique will render the process of protein structure determination faster and
109 easier. Haruspex was trained on a carefully curated ground truth dataset based on experimental
110 data from EMDB. The pre-trained network can be straightforwardly applied to predict structures in
111 newly reconstructed Cryo-EM density maps, and will be refined and adapted as new data become
112 available. Besides guidance for model building and domain placements, the network may also be
113 useful for model validation due to its high median recall and precision rates of 94.2% and 79.2%,
114 respectively. The trained network and documentation are available from
115 gitlab.com/phimos/haruspex.

116 **Acknowledgements**

117 This work was supported by the DFG (project TH2135/2-1), the High Performance Computing Cloud
118 of Würzburg University (DFG project 327497565) and the Rudolf Virchow Center for Experimental
119 Biomedicine. We would like to thank Bettina Böttcher, Niko Grigorieff, Tom Burnley and Jola Mirecka
120 for fruitful discussions; and Bernhard Fröhlich for great computational support.

121 *Methods*

122 *Training Data*

123 We queried the Electron Microscopy Data Bank (EMDB) for all single particle Cryo-EM maps with a
124 resolution ≤ 4 Å, for which corresponding protein models were available in the Protein Data Bank in
125 Europe (PDBe), yielding 576 map and model pairs as of February 2018. We filtered these EMD/PDB
126 pairs by the following three criteria: (1) Good fit between map and model; (2) presence of at least
127 one annotated α -helix or β -sheet; and (3) preference of the highest resolution maps in case the
128 same authors deposited several instances of the same macromolecular complex. Maps with severe
129 misfits, misalignments, or models without corresponding reconstruction densities, and vice versa,
130 were discarded. After applying these criteria, we retained 293 map/model pairs for generating the
131 training data.

132 To extract secondary structure information from the PDB data, we developed a custom parser for
133 the PDBML [19] format based on xmldict [20]. To obtain additional secondary structure
134 information, we implemented a variant of the DSSP algorithm [11] without strand direction, and a
135 torsion angle based secondary structure detection inspired by STRIDE [13]: annotated or DSSP-
136 detected secondary structures were extended by neighbouring amino acids if they matched the
137 same Ramachandran profile.

138 *Annotation of reconstruction maps*

139 For every entry pair, the augmented model was then superimposed on the map and all voxels within
140 3 Å of a C α or C,N,O-backbone atom, or, in the case of nucleotides, within 3 Å of any atom, were
141 assigned the respective class (helix, sheet or nucleotide) if their value was higher than $\frac{1}{2}$ of the
142 average backbone density of the helix, sheet or nucleotide in question. Secondary structures with a
143 backbone standard deviation of $< 2 \sigma$ and atoms without secondary structure assignment were
144 labelled as “empty” to exclude incorrectly modelled, misfitted, or flexible structures. For some
145 training data pairs (e.g. virus capsids), only small or partial protein models were deposited for large
146 Cryo-EM maps, resulting in well-defined high-density regions without model coverage. These regions
147 will not get annotated and will result in false positives if the network tries to predict the actual
148 structure. To mitigate this, all voxels with density $\geq 1.0 \sigma$ but not within 5 Å of a model atom with
149 density $\geq 1.0 \sigma$ were masked as unmodeled density and hence did not contribute to training.

150 Since our network generated a single class label as output, the reconstruction density of the
151 secondary structures must be converted to a strict assignment to one of the three classes in order to
152 be used as training examples. For each secondary structure, the reconstruction map density was
153 multiplied by the backbone standard deviation and rescaled to an output density between zero and
154 one (corresponding to 0.5 and 1.0 times the average backbone density of the local secondary
155 structure element) for each label type. The highest channel value determined the voxel class. If
156 multiple channels shared the same value, sheets took precedence over nucleotides, which took
157 precedence over helices. Voxels where all channel values were below 0.01 were assigned the
158 “empty” class. Finally, reconstruction maps were rescaled to a voxel size of 1.1 Å, if they were
159 outside of [1.0; 1.2] Å.

160 *Generation of training segments*

161 To generate the 70^3 voxel sized segments needed for training, candidate volumes were sampled
162 from the entire map, and segments with a mean backbone density $<3.0 \sigma$, less than 5% annotated
163 volume, or less than 100 atoms with standard deviation $\geq 1.0 \sigma$ were discarded. This resulted in
164 altogether 2183 training segments, of which 110 segments (5%) were held back for evaluation
165 during training. To generate additional segments for training, we applied rotations in steps of 90°
166 around all three axes, resulting in 24 rotated versions of each segment that could all be used as
167 separate training volumes since the convolutional network is not rotation-invariant. Segments were
168 further augmented during training by using a randomly translated 40^3 sub-cube for each step.

169 *Network Architecture*

170 We use a 3D U-Net architecture with a single input channel (reconstruction density) and an input
171 layer size of 40^3 voxels, shown in supplementary Fig. 1. The encoding branch consisted of two $3 \times 3 \times 3$
172 convolutional layers with 32 and 64 feature channels, respectively, followed by $2 \times 2 \times 2$ max-pooling
173 layers. Another convolutional layer with 128 feature channels followed by $2 \times 2 \times 2$ max-pooling layer
174 finally resulted in an 8^3 cube with 128 feature channels at the deepest layer of the network. This
175 cube was passed through another convolutional layer with the same data padding in order to
176 preserve its dimensions. A fully connected layer was considered, but not chosen due to its high
177 memory and performance cost. The decoding branch of the U-Net was made of two blocks, each
178 consisting of a deconvolution followed by two $3 \times 3 \times 3$ convolutions (128 feature channels in the first,
179 64 and 32 channels in the second block to restore symmetry) with concatenated sections of the
180 corresponding layer in the encoding part. The output part consists of a final $1 \times 1 \times 1$ convolution
181 followed by a soft-max output layer. The output layer reproduced the central 20^3 voxel cube of the
182 input layer in four annotation channels representing co-dependent probabilities for the four classes
183 (helix, sheet, nucleotide, empty) summing up to one. The highest channel value determined the
184 predicted class. Implementation was realized using Tensor Flow [21].

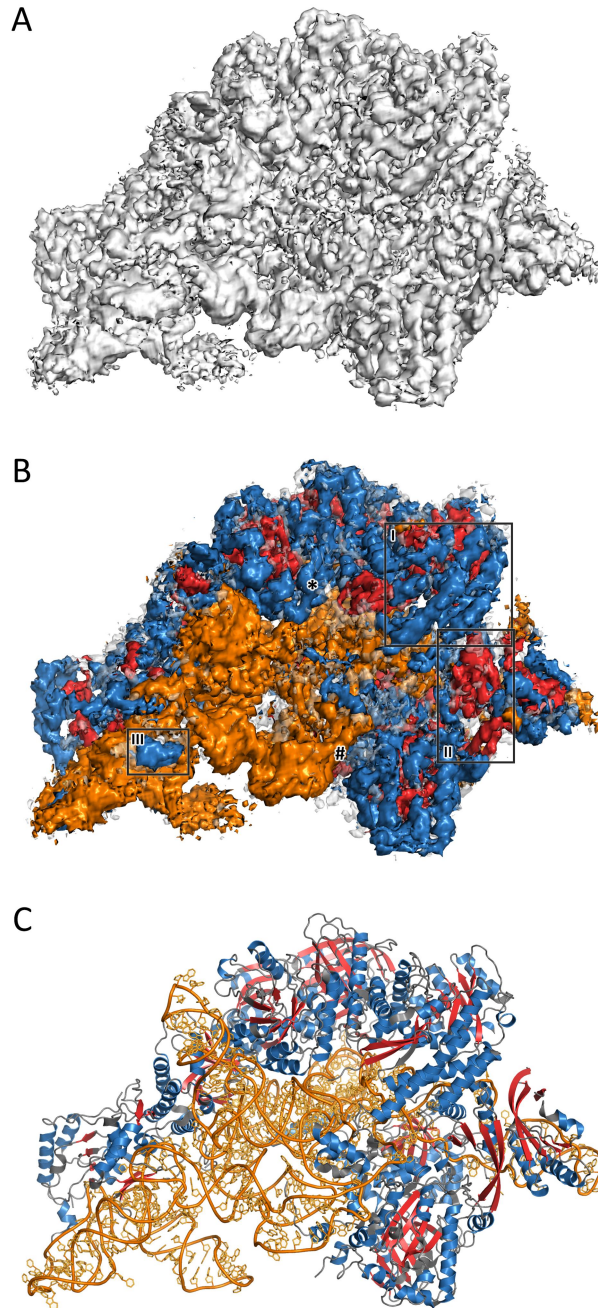
185 *Network Training*

186 The network was trained for 40,000 steps on training batches of 100 random segment pairs per step,
187 using ADAM stochastic optimization [22] with a learning rate of 0.001, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\varepsilon =$
188 0.1. Error assignment for backpropagation was performed using cross-entropy loss, where the target
189 class was represented in one-hot encoded binary format (1 for the target class, 0 for the other three
190 classes). To account for class imbalance, voxels were weighted according to overall class occurrence
191 in the training data. Furthermore non-true negatives were weighted 16-fold stronger than true
192 negatives due to an overabundance of the latter.

193 **References:**

- 194 [1] Cireşan, D. C., Meier, U., Masci, J. & Schmidhuber, J. (2011). *The 2011 International Joint*
195 *Conference on Neural Networks*. 1918–1921.
- 196 [2] Cireşan, D. C., Giusti, A., Gambardella, L. M. & Schmidhuber, J. (2013). *Medical Image Computing*
197 *and Computer-Assisted Intervention – MICCAI 2013*, Vol. pp. 411–418. Springer, Berlin, Heidelberg.
- 198 [3] Çiçek, Ö., Abdulkadir, A., Lienkamp, S. S., Brox, T. & Ronneberger, O. (2016). *International*
199 *Conference on Medical Image Computing and Computer-Assisted Intervention*, Vol. pp. 424–432.
200 Springer.
- 201 [4] Falk, T., Mai, D., Bensch, R., Çiçek, Ö., Abdulkadir, A., Marrakchi, Y., Böhm, A., Deubner, J., Jäckel,
202 Z., Seiwald, K., Dovzhenko, A., Tietz, O., Dal Bosco, C., Walsh, S., Saltukoglu, D., Tay, T., Prinz, M.,
203 Palme, K., Simons, M., Diester, I., Brox, T. & Ronneberger, O. (2019). *Nature Methods*. **16**,
- 204 [5] Jiang, W., Baker, M. L., Ludtke, S. J. & Chiu, W. (2001). *J Mol Biol*. **308**, 1033–1044.
- 205 [6] Kong, Y., Zhang, X., Baker, T. S. & Ma, J. (2004). *Journal of Molecular Biology*. **339**, 117.
- 206 [7] Baker, M. L., Ju, T. & Chiu, W. (2007). *Structure*. **15**, 7.
- 207 [8] Shelhamer, E., Long, J. & Darrell, T. (2017). *IEEE Transactions on Pattern Analysis and Machine*
208 *Intelligence*. **39**, 640–651.
- 209 [9] Tagari, M., Newman, R., Chagoyen, M., Carazo, J.-M. & Henrick, K. (2002). *Trends in Biochemical*
210 *Sciences*. **27**, 589.
- 211 [10] Berman, H., Henrick, K. & Nakamura, H. (2003). *NSMB*. **10**, 980.
- 212 [11] Kabsch, W. & Sander, C. (1983). *Biopolymers: Original Research on Biomolecules*. **22**, 2577–
213 2637.
- 214 [12] Frishman, D. & Argos, P. (1995). *Proteins: Structure, Function, and Bioinformatics*. **23**, 566–579.
- 215 [13] Tronrud, D. E., Berkholtz, D. S. & Karplus, P. A. (2010). *Acta Crystallographica Section D Biological*
216 *Crystallography*. **66**, 834–842.
- 217 [14] Hollingsworth, S. A. & Karplus, P. A. (2010). *BioMolecular Concepts*. **1**,
- 218 [15] Si, D., Ji, S., Nasr, K. A. & He, J. (2012). *Biopolymers*. **97**, 698–708.
- 219 [16] Li, R., Si, D., Zeng, T., Ji, S. & He, J. (2016). 2016 IEEE International Conference on Bioinformatics
220 and Biomedicine (BIBM), Vol. pp. 41–46.
- 221 [17] Tang, G., Peng, L., Baldwin, P. R., Mann, D. S., Jiang, W., Rees, I. & Ludtke, S. J. (2007). *Journal of*
222 *Structural Biology*. **157**, 38–46.

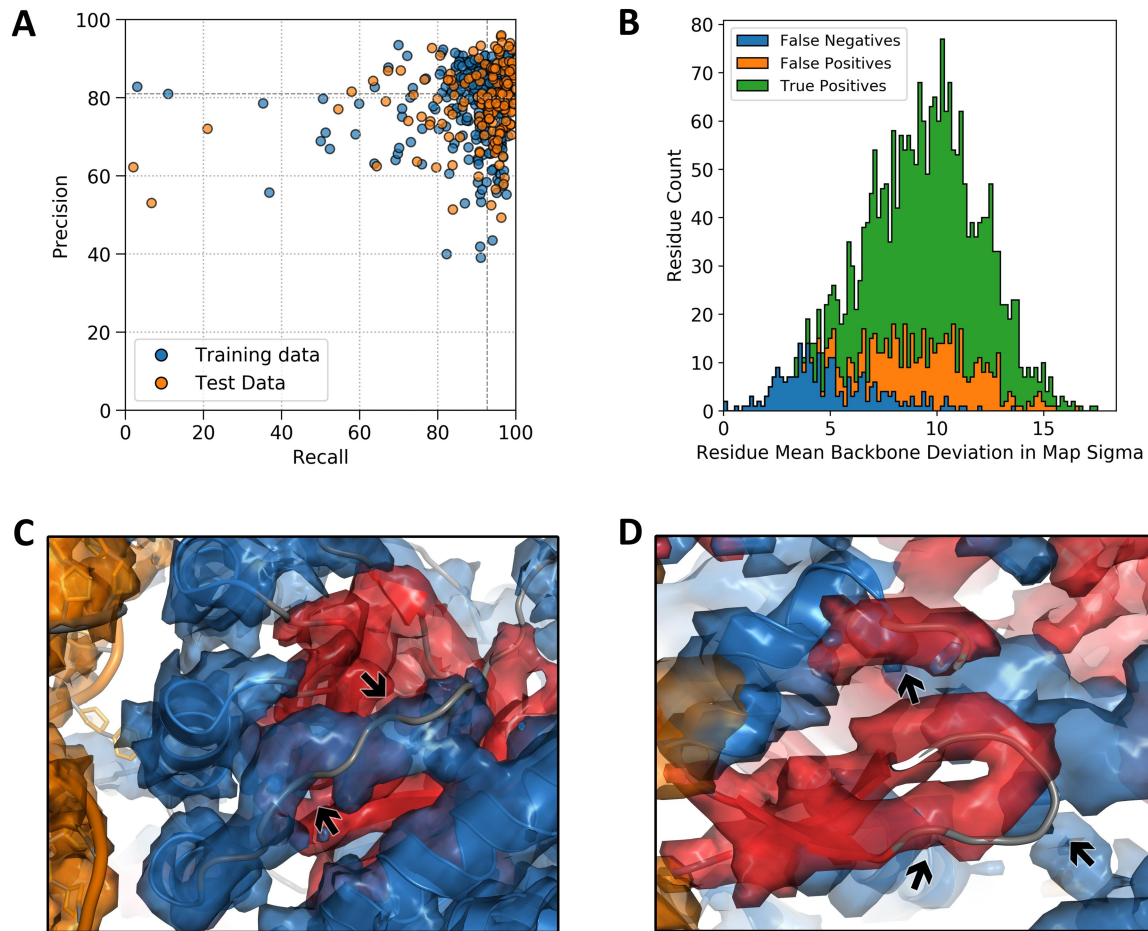
- 223 [18] Haslam, D., Zeng, T., Li, R. & He, J. (2018). Proceedings of the 2018 ACM International
224 Conference on Bioinformatics, Computational Biology, and Health Informatics, Vol. pp. 628–632.
225 New York, NY, USA: ACM.
- 226 [19] Westbrook, J., Ito, N., Nakamura, H., Henrick, K. & Berman, H. M. (2004). *Bioinformatics*. **21**,
227 988–992.
- 228 [20] Blech, M. xmltodict.
- 229 [21] Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G., Davis, A., Dean, J.,
230 Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser,
231 L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M.,
232 Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F.,
233 Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y. & Zheng, X. (2015). TensorFlow: Large-
234 Scale Machine Learning on Heterogeneous Distributed Systems.
- 235 [22] Kingma, D. P. & Ba, J. (2014). *ArXiv Preprint ArXiv:1412.6980*.
- 236



237

238 **Figure 1. Haruspex Annotation.** **A.** Reconstruction map for the human Ribonuclease P holoenzyme
239 (EMDB entry 9627). Manual assignment of secondary structure features can be difficult, in particular
240 if the composition of a macromolecular complex is unknown. The shown surface corresponds to $\sigma =$
241 0.04 with no carving. **B.** Secondary structure as identified by our network in the map, projected onto
242 the surface. Orange corresponds to nucleotides; blue to helices; red to sheets and transparent grey
243 were not assigned any secondary structure. This was a fairly typical test case with 70.5% true
244 positives, 18.8% false positives and 10.7% false negatives. Recall was 86.8% and precision 79.0%.
245 Region (I) depicts a well-predicted α -helical structure, (II) a β -sheet and (III) RNA misinterpreted as

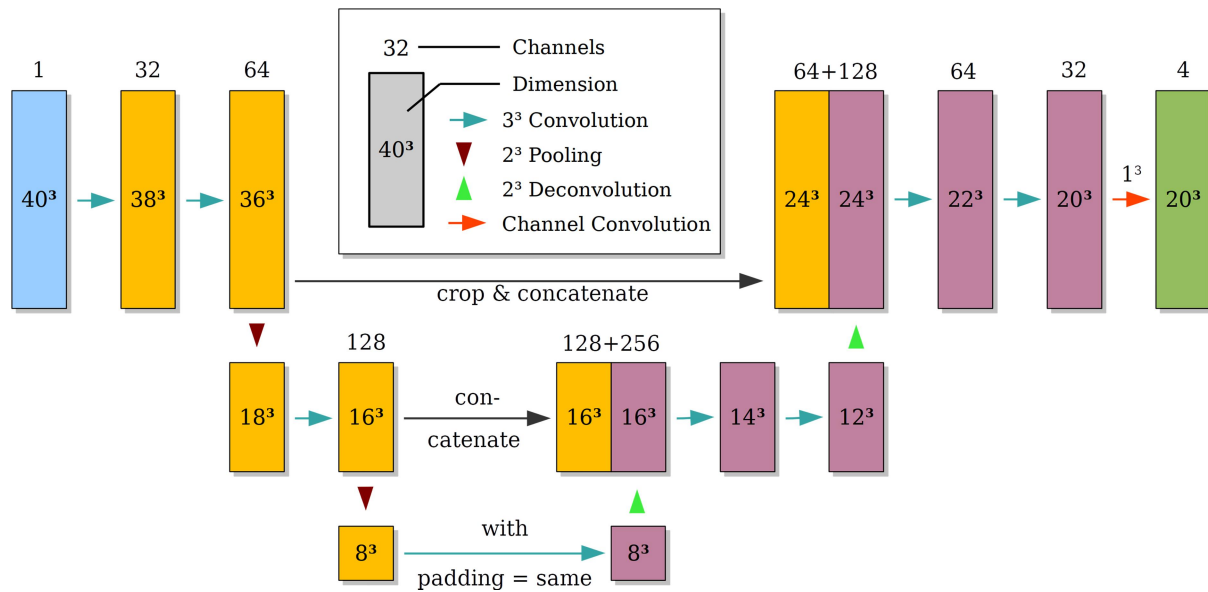
246 an α -helix. **C.** Model (PDB 6AHU) for comparison. The regions depicted in Fig. 2C and 2D are marked
247 # and *, respectively.



248

249 **Figure 2. Network performance.** **A.** Network precision vs. recall rates, with one marker per EMDB
250 entry (test set markers are orange, training set markers blue). Both perform similarly well. **B.**
251 Frequency vs. map σ level for EMDB 9627 on a per-residue basis: True positives (green), false
252 positives (orange) and false negatives (blue). This plot is typical: false negatives often occur in low
253 density map regions. **C.** α -Helical false positives (PDB 6AHU, J131 – 139): The model partly occupies
254 the conformational space of a polyproline type II helix (P_{II}), which is often misinterpreted as α -helical
255 and may have been modelled incorrectly (given that the model does not completely fit the density).
256 **D.** False positives in a β -sheet (6AHU, B215-B221). The deposited model does not maintain the
257 hydrogen bonding that defines β -sheets; to the network, however, the fold still 'looks' like a β -sheet
258 and a third segment (top) is also assumed to be part of it.

259



260
261

262 **Figure 3 (Methods). Haruspex neural network architecture.** The network consists of multiple
 263 interconnected layers, shown as rectangular boxes. We employed a state-of-the-art U-Net-like
 264 encoder-decoder architecture [4], a subclass of so-called fully-connected networks where spatial
 265 information and object details are encoded, reduced by pooling layers and then recovered again
 266 with up-sampling or transpose convolutions. The term U-Net arises from the U-like shape of the data
 267 flow. The layers are connected by convolution and pooling operations (arrows). Layer height
 268 represents the level of abstraction: lower layer data, generated by pooling operations, contain more
 269 abstract representations of the map. Input data (blue) is fed into the downconvolutional arm
 270 (yellow) in order to extract valuable information, which is then combined with previously discarded
 271 information through concatenations in the upconvolutional arm (purple) to compute annotated
 272 output data (green) for a subsection (20^3) of the input volume (40^3). Our network consists of two
 273 encoder blocks, containing altogether three convolutional layers ($3 \times 3 \times 3$) and two pooling layers
 274 ($2 \times 2 \times 2$). This is followed by two decoder blocks, one with upconvolution followed by two $3 \times 3 \times 3$
 275 convolutions and 128 feature channels, and one with upconvolution followed by two $3 \times 3 \times 3$
 276 convolutions with 64 and 32 feature channels, with concatenated sections of the corresponding
 277 layer in the encoding part. The output part consists of a final $1 \times 1 \times 1$ convolution followed by a soft-
 278 max output layer. This results in 13 layers in total (12 + 1 convolution at bottom). The network is
 279 trained end-to-end by comparing the predicted class of each voxel to the annotated EMDB model
 280 using cross-entropy loss, propagating the error back through the network, and adapting the network
 281 weights to iteratively minimize the error.