# Model building and validation for cryoEM

*Oliver Clarke*

*(@OliBClarke)*

COLUMBIA UNIVERSITY
MEDICAL CENTER

**"Is my map buildable??"** 🤷‍♂️

# General principles to keep in mind…

- Be skeptical of your own data – always consider the null hypothesis
  - E.g. particles **not** protein, not **my** protein, ligand **not** bound, etc…

- Test, test, & test again…
  - E.g. do different reconstruction/refinement algorithms give same result?
  - Confirm presence of ligands/subunits using difference maps, labels, mass spec..

- Trust your eyes!
  - If 3Å reconstruction has no sidechains, retrace your steps – postprocessing? Masking? Initial model?
  - Remember that many metrics (e.g. FSC) measure **self-consistency**, not accuracy/reality. No substitute (yet!) for human judgement and careful attention to detail.

**An atomic model is a compact interpretation of the density map in light of prior knowledge (both specific and general).**

- Aim is to build a model that is consistent with **both** the density map and everything we independently know about the structure & composition of the macromolecule of interest, both specifically and in terms of our general knowledge of protein structure and chemistry.

- At medium resolution (3-5 Å), this still requires manual building (yes, even if you start from an AlphaFold prediction...🙃). Even the best autobuilt model still requires manual inspection in **all**, and correction in *most* cases. (generates many fragments which need inspection, correction, merging)

- Tradeoff between available prior knowledge and required resolution for atomic modelling – at the extremes, if a complete crystal structure is already available, 10Å data may be sufficient, while if no sequence/composition data is available even 3Å may not suffice.

**Prior knowledge**
- Protein sequence and derived info (secondary structure predictions, covariation/conservation, patterns of large/aromatic residues), disorder & contact prediction
- Crystal structures (+ homology & **ML-derived models – AlphaFold, Modelangelo**)
- Knowledge of protein structure, folding, chemistry, geometry.

**Density map**
- Resolution (+ local resolution, + map modification/sharpening)
- Patterns of large/small/absent sidechains
- Sharpening and density modification
- Conformational/compositional heterogeneity

**Building & refinement
(Chimera, COOT, ISOLDE, etc…)**

**Atomic model**

- If possible, unique model (or ensemble) that agrees with both density map and priors
- Otherwise (and per region), specify ambiguity (w/UNK residues and numbering or Ca only model)
- Validation not just (or even mostly) about overfitting.
- Identify, analyse, fix errors.
- Direction and register of sequence fit.
- Ligand identification/assignment.
- No model is, or ever will be perfect. That's okay.

**One extreme – at atomic resolution, the position of many atoms can be inferred without prior knowledge of the sequence**

Yip, K.M., Fischer, N., Paknia, E. et al. Atomic-resolution protein structure determination by cryo-EM. Nature **587,** 157–161 (2020)

**At 20 Å (here using cryoET), an informative model can be generated by taking advantage of external information – crystal structures, connectivity from crosslinking & MS, even when de novo building is not possible.**

Kim, S., Fernandez-Martinez, J., Nudelman, I. et al. Integrative structure and functional anatomy of a nuclear pore complex. Nature **555,** 475–482 (2018)

Usually, we are somewhere in between the two – combining prior knowledge with inferences made from analyzing the density map.

To build a better/more reliable model, we can either get additional/better priors, or improve our density map (or part of it).

**Before you start – make sure your maps are appropriately sharpened and low pass filtered! (and consider whether building is justified or whether further improvement of the reconstruction is required first)**

- **Often it is helpful to build using multiple maps**. Assuming 3-3.5Å global res, I would suggest using a map filtered to the global resolution, one filtered to the best local resolution, and one filtered to ~4-4.5 Å (to better visualize connectivity and mobile ligands/lipids).

- Try both simple B-factor sharpening and the approach used by `phenix.resolve_cryo_em`, which incorporates anisotropy removal and statistical density modification. In cases of **severe** anisotropy, deepEMhancer/EMReady/spIsoNet can be useful to assist map interpretation (**approach with caution**).

- Also, if your map doesn't "look like" 4 Å, trust your eyes! If it is nominally 4Å and there are no sidechains visible, or your helices look "stretched", assess orientation bias (3D-FSC server: *https://3dfsc.salk.edu*), local resolution variation, and double check sharpening and masking parameters (are you *sure* you're looking at the sharpened map? Is the mask used for FSC calculation sensible?)

**Example of map anisotropy mitigated by masked refinement**

- Map anisotropy hinders interpretation, even when resolution in "good" direction is high

- Can derive from either preferred orientation, or interdomain mobility (or combination).

- In latter case, masked refinement can improve local map quality to aid model building and map interpretation. **Always better to improve the map than build in marginal density**

- If anisotropy derives from preferred orientation, it is best to address this by improving the sample (additives, constructs, substrates) or altering data collection strategy (tilt). If all else fails, ML-based map improvement using deepEMhancer can improve map interpretability.

# Map anisotropy can be mitigated by masked refinement

- Map anisotropy hinders interpretation, even when resolution in "good" direction is high

- Can derive from either preferred orientation, or interdomain mobility (or combination).

- In latter case, masked refinement can improve local map quality to aid model building and map interpretation. Always better to improve the map than build in marginal density

- If anisotropy derives from preferred orientation, it is best to address this by improving the sample or data collection (tilt). **If all else fails, ML-based map improvement using deepEMhancer, spIsoNet or EMReady can improve map interpretability (caveats abound!).**
- Also check out new approach from Yifan Cheng's lab – non-ML based map deconvolution, AR-Decon



DeepEMhancer

*(Sanchez-Garcia R., 2021 Comm. Biol.)*

# Example from the literature



**Good direction – looks okay**

# Example from the literature



**Bad direction – "ropy", stretched, hard to interpret**

**Prep for model building - what can we learn from the sequence alone?**

Your protein sequence contains a lot of useful information which you can use to aid model building:

- Start by identifying boundaries of conserved domains (NCBI CDD: https://www.ncbi.nlm.nih.gov/Structure/cdd/; DELTA-BLAST also performs CD-search by default)

- Then identify and/or generate suitable structural templates for building known domains: Predominantly alphafold, unless a better experimentally determined model is available. Modelangelo also useful for generating initial model, especially if sequence unknown.

- Secondary structure, TM & disorder prediction (XtalPRED for overall summary; specific tools such as SPOT-DISORDER, SPIDER3 for best accuracy).

- Contact prediction from evolutionary couplings: EVFOLD & GREMLIN.

- Conservation analysis: Use favorite MSA algorithm (MUSCLE & CLUSTAL-OMEGA work well; TM-COFFEE, PRALINE-TM useful for membrane proteins) to create a sequence alignment of your protein with a few orthologs; gaps & insertions most commonly occur in loops/disordered regions. Useful as a guide during building.

**Prep for model building - what can we learn from the sequence alone?**

Your protein sequence contains a lot of useful information which you can use to aid model building:

- **Start by identifying boundaries of conserved domains (NCBI CDD: https://www.ncbi.nlm.nih.gov/Structure/cdd/; DELTA-BLAST also performs CD-search by default)**

- Then identify and/or generate suitable structural templates for building known domains: FUGUE, PHYRE2, MUSTER. (Alphafold/ROSETTAfold dominate now!). Modelangelo useful for generating initial model, especially if sequence unknown.

- Secondary structure, TM & disorder prediction (XtalPRED for overall summary; specific tools such as SPOT-DISORDER, SPIDER3 for best accuracy).

- Contact prediction from evolutionary couplings: EVFOLD & GREMLIN.

- Conservation analysis: Use favorite MSA algorithm (MUSCLE & CLUSTAL-OMEGA work well; TM-COFFEE, PRALINE-TM useful for membrane proteins) to create a sequence alignment of your protein with a few orthologs; gaps & insertions most commonly occur in loops/disordered regions. Useful as a guide during building.

**CDD provides a guide to domain level architecture, including sequence alignments & representative structures.**



TED (The Encyclopedia of Domains) also a good option for this.

**CDD provides a guide to domain level architecture, including sequence alignments & representative structures.**
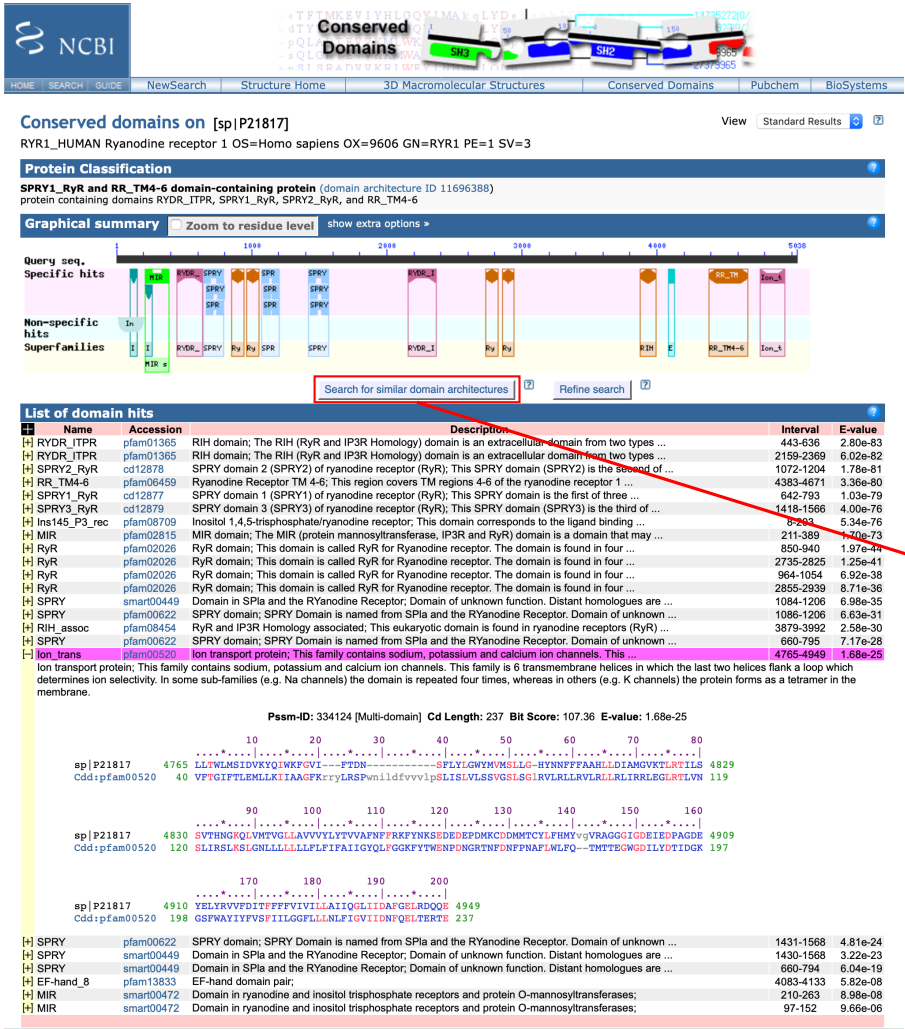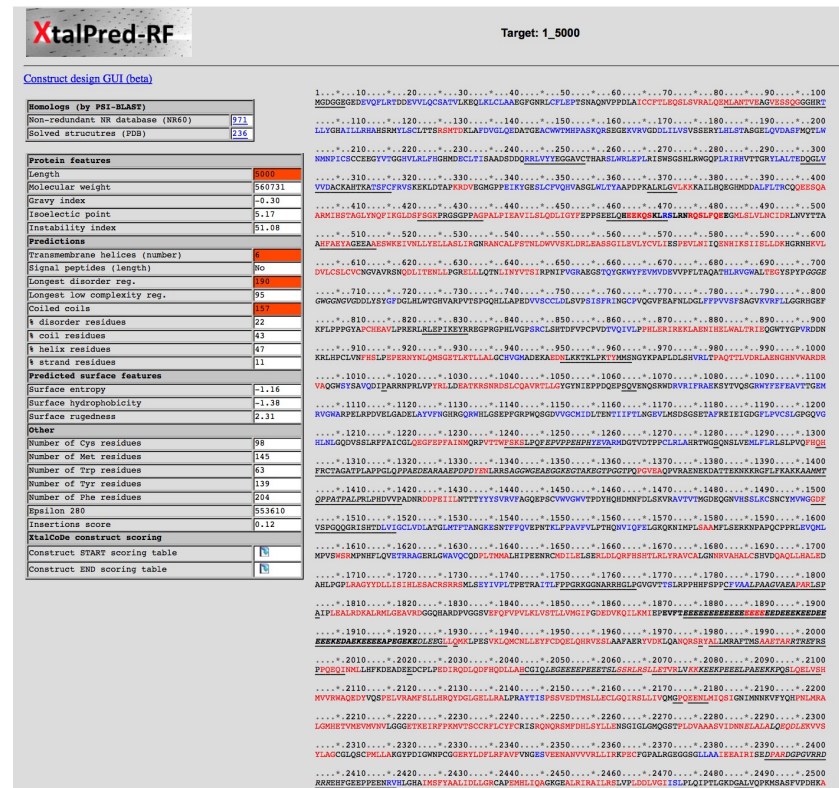
**CDD provides a guide to domain level architecture, including sequence alignments & representative structures.**



Once an initial trace is obtained for these regions, use Foldseek to identify structural homologs that could not be identified by sequence alone.

**CDD provides a guide to domain level architecture, including sequence alignments & representative structures.**



Useful for identifying related proteins – which one can then investigate further using Alphafold.

# XtalPRED is a great tool for summarizing predicted sequence properties.



Highlights predicted secondary structure, disorder, low complexity regions on sequence in an easily digestible format. Useful to print and consult while building. Also provides list of structural homologs. (**http://ffas.burnham.org/XtalPred-cgi/xtal.pl**)

# Secondary structure prediction is a very useful guide when building.





**Where is this motif in the sequence?**

# Secondary structure prediction is a very useful guide when building.



Secondary structure prediction is ~80% accurate. So if your model consistently disagrees with predicted secondary structure, look at it very closely!

**What can we learn from the map alone?**

**What can we learn from the map alone?**



**Left handed! Obvious here – can be less clear at lower res, so be careful.**

**OK, that's better! What can we learn from the map alone?**

# Which direction does the helix point?

**Which direction does the helix point?**

# Helices – alpha and $3_{10}$



**Alpha**

- ~90%
- 3.6 residues per turn
- Fat

**$3_{10}$**

- ~10%. More common in TM? (e.g. S4 of VSD)
- 3 residues per turn. Triangular cross section.
- Skinny
- Can be tricky to identify at low resolution, can lead to register errors.

**Can we identify any probable sidechains from the density?**

**Can we identify any probable sidechains from the density?**

**Can we identify any probable sidechains from the density?**

**Can we identify any probable sidechains from the density?**

**Test the initial hypothesis by extending sequence assignment along the chain.**



...VFNSLTEYIQGPCTGNQQSLAHSRL**W**DAVVG**F**LHV**F**AHMMMKLAQDSSQIELLKELLDLQ...

**Test the initial hypothesis by extending sequence assignment along the chain.**



...VFNSLTEYIQGPCTGNQQSLAHSRL**WDAVVGFLHVFAHMMMK**LAQDSSQIELLKELLDLQ...

**Test the initial hypothesis by extending sequence assignment along the chain.**

Notice that the **absence** of large sidechain densities at small residue positions is just as valuable in validating the fit as the fit of large sidechains to the density.



...VFNSLTEYIQGPCTGNQQSLAHSRL**WDAVVGFLHVFAHMMMK**LAQDSSQIELLKELLDLQ...

**Test the initial hypothesis by extending sequence assignment along the chain.**

Also, note that the information content of local regions varies. Consider "VTVVAASSTVV" vs "FGAAYWVTRA" – which is more likely to be uniquely identifiable from the map?



...VFNSLTEYIQGPCTGNQQSLAHSRL**WDAVVGFLHVFAHMMMK**LAQDSSQIELLKELLDLQ...

**CryoID can help when you don't even know the sequence!**

- Similar approach codified and automated in the "cryoID" program – but in this case, starting from the density, with no sequence input!

- Split map into fragments

- Use reduced complexity pseudo-sequence to convert map fragments into motifs which can be used to search sequence database.

- Identify most likely candidate sequence, combine fragments and rebuild.

- Useful when purifying from endogenous sources, where composition may not be known.



*(Ho et al., Nature Methods, 2020)*

**What if we don't have sidechains?**

- CryoID requires sidechains – what if we don't have them? E.g. sub-nm reconstructions from cryoET

- Still a lot of information encoded in arrangement of secondary structural elements.

- Can use COLORES (in SITUS package) to query Alphafold database using segmented densities as query

- Requires caution and biochemical validation in interpreting results (and subsequent fitting).



Figure 2. Unbiased matching of target densities from *in situ* cryoET reconstructions with structure models of 21,615 mouse proteins predicted using AlphaFold2.

*(Zhen et al., BiorXiV, 2023)*

# Modelangelo – applying neural networks to map interpretation

Here, we introduce a machine-learning approach, called ModelAngelo, for the automated building of atomic models and the identification of proteins in cryo-EM maps. Machine learning approaches often require large amounts of training data. For example, recent protein language models were trained on tens of millions of sequences (*14*) and AlphaFold2 was trained on more than 200,000 structures (*15*). In contrast, fewer than 13,000 cryo-EM structures with resolutions better than 4 Å have been determined to date and many of these are redundant. The limited amount of available training data prompted us to design a multi-modal machine-learning approach that combines local information from the cryo-EM map surrounding each protein or nucleic acid residue with additional information from the protein sequences in the sample and the local geometry of the structure. Similar sources of information are exploited by human experts when manually building atomic models in cryo-EM maps.

*(Jamali et al., BiorXiv, 2023)*

**Modelangelo – applying neural networks to map interpretation**



Step 1
Predicted residue positions

Step 2
Initialize graph with
nearest neighbours

Optimize graph with
Graph Neural Network

Step 3
Postprocessed model

*(Jamali et al., BiorXiv, 2023)*

# Modelangelo – applying neural networks to map interpretation



*(Jamali et al., BiorXiv, 2023)*

# Modelangelo – applying neural networks to map interpretation



*(Jamali et al., BiorXiv, 2023)*

# Modelangelo – applying neural networks to map interpretation



*(Jamali et al., BiorXiv, 2023)*

**Modelangelo – applying neural networks to map interpretation**

- ***Fantastic*** starting point for model building – generates near-complete models with good geometry in favorable areas.

- Difficult regions (e.g. flexible, anisotropic) still require manual building (for now!)

- Model still requires manual inspection and analysis (both for completion/correction/validation, and understanding!)

- Can't build waters/ions/ligands (yet!)

- ***Allows us to build better, faster, and focus on difficult/important regions.***

*(Jamali et al., BiorXiv, 2023)*

**How to deal with uncertainty in sequence assignment and sidechain placement**

- You will likely encounter situations where you cannot be certain of the local sequence register – what to do?

- No clear consensus, but I suggest assigning residue code as "UNK" and numbering to "best guess" value. A more granular way to quantify/convey uncertainty would be helpful!

- Sidechain placement – two main camps – trim sidechains to density vs place them all (+/- zero occ.). The former may sound more conservative, but it can hide errors during validation (during analysis of clashes).

- Either is acceptable, just be consistent, and preferably document the approach taken when writing up the structure. You can read long discussions of the merits of both approaches on the CCP4bb should you so desire 🙃

**Prior knowledge can come in many forms – use any and all available info to guide model building.**



Here, serendipitous identification of a conformational class of RyR1 lacking density for one subunit aided identification of protomer boundaries. In other cases, cross-linking data or NS data on subcomplexes or Fab-complexes may be helpful.

*(Zalk et al, Nature 2015)*

**In a similar manner, we can use locally aligned difference maps between holo and apo structures to locate ligands.**

# The three ligands are clustered around the C-terminal domain.



(Ca$^{2+}$ only) minus (EGTA only)

# The three ligands are clustered around the C-terminal domain.



**(ATP/Caffeine) minus (EGTA only)**

**EM-specific considerations**

- No unambiguous sequence/elemental markers at low resolution (no equivalent of SeMet yet).

- No feedback from phase improvement, but also no model bias – WYSIWIG.

- Often substantial variation in local resolution – different strategies and levels of detail required for different regions. Map sharpening essential.

- "Medium" resolution (4-6Å) much more common than for crystallography.

- Often have more than one map, with different composition or conformation (may be convenient to combine focused refinements in Chimera by taking max value at each voxel after alignment, e.g.: `vop maximum #1,2 ongrid #1` )

**Building an initial model - where to start?**

- If you have a crystal structure, of a fragment or a homology model of a domain, place it, and extend into density. (***Now, Alphafold & Rosettafold mean this is almost always the case***)

- Otherwise, identify structurally distinctive motifs in the sequence – for example, a strongly predicted helix with three aromatic residues near the N-term end – and identify candidate locations in the density map. Extend and see if hypothesis still holds.

**Using UCSF Chimera to fit solved domains**



**Start with map and model.**

**Using UCSF Chimera to fit solved domains**



**Move model to approximate position (if known, to save computation)**

**Using UCSF Chimera to fit solved domains**



**Run fitmap with 'search' (here 100 orientations) and 'radius' (here 5 Å)**

## Using UCSF Chimera to fit solved domains



**Chimera will return a list of candidate orientations, ranked by agreement with the map. Hopefully there will be a clear separation between the correct and incorrect solutions.**

# Using UCSF Chimera to fit solved domains



Chimera will return a list of candidate orientations, ranked by agreement with the map. Hopefully there will be a clear separation between the correct and incorrect solutions.

**Using UCSF Chimera for voxel size calibration (of your map and others)**

- Voxel size generally requires calibration against a crystal structure.

- Once calibrated, generally stable between samples/datasets at same magnification.

- Can calibrate by fitting in Chimera at range of nominal voxel sizes and measuring correlation.

- Incorrect voxel sizes are common in deposited maps - **be aware of this when comparing structures**. E.g. here there is a 3% difference – affects structural alignment, reported resolution (3.8 vs 3.9Å).



**1.34 Å (original)**
**1.38 Å (corrected)**

**COOT – Crystallographic Object Oriented Toolkit**

- Simple, intuitive interface for building and manipulating atomic models in density maps.

- Low computational requirements

- Extensive API – easy to script or modify (using simple Python code)

- On-the-fly sharpening, resampling and low pass filtering



(Try the latest nightly with new features for EM, improved RSR: http://www.ccpem.ac.uk/download.php)

(*Emsley P. 2004, Acta Cryst. D; Casañal A. et al. 2020, Protein Science*)

## COOT – Crystallographic Object Oriented Toolkit

- Simple, intuitive interface for building and manipulating atomic models in density maps.

- Low computational requirements

- **Extensive API – easy to script or modify (using simple Python code)**

- On-the-fly sharpening, resampling and low pass filtering

**COOT – Crystallographic Object Oriented Toolkit**

- Simple, intuitive interface for building and manipulating atomic models in density maps.

- Low computational requirements

- **Extensive API – easy to script or modify (using simple Python code)**

- On-the-fly sharpening, resampling and low pass filtering

**Sidenote – trimmings for ChimeraX:**

- Added some shortcuts and a bunch of aliases that I find helpful for using ChimeraX with maps/models.
- Use and modify as you like!

**Sidenote – trimmings for ChimeraX:**

**COOT – Crystallographic Object Oriented Toolkit**

- Simple, intuitive interface for building and manipulating atomic models in density maps.

- Low computational requirements

- **Extensive API – easy to script or modify (using simple Python code)**

- On-the-fly sharpening, resampling and low pass filtering

```python
def mutate_by_entered_code():
  def mutate_single_letter(X):
    entry=str(X).upper()
    mol_id=active_residue()[0]
    ch_id=active_residue()[1]
    resno=active_residue()[2]
    ins_code=active_residue()[3]
    resname=residue_name(mol_id,ch_id,resno,ins_code)
    map_id=imol_refinement_map()
    aa_dic={'A':'ALA','R':'ARG','N':'ASN','D':'ASP','C':'CYS','E':'GLU','Q':'GLN','G':'GLY','H':'HIS','I':'ILE','L':'LEU','K':'LYS
    nt_list=['A','C','T','G','U']
    if (resname in aa_dic.values()) and (aa_dic.get(entry,0)!=0):
      mutate(mol_id,ch_id,resno,ins_code,aa_dic.get(entry,0))
    elif (resname in nt_list) and (entry in nt_list):
      mutate_base(mol_id,ch_id,resno,ins_code,entry)
    else:
      info_dialog("Invalid target residue! Must be protein or nucleic acid, and entered code must be single letter.")
  generic_single_entry("New residue? (single letter code)","A","Mutate by single-letter code",mutate_single_letter)
```

```python
#mutate active residue to entered residue code (upper or lower case single-letter)
add_key_binding("Mutate by single letter code","M",
lambda: mutate_by_entered_code())
```

**COOT – Crystallographic Object Oriented Toolkit**

- Simple, intuitive interface for building and manipulating atomic models in density maps.

- Low computational requirements

- **Extensive API – easy to script or modify (using simple Python code)**

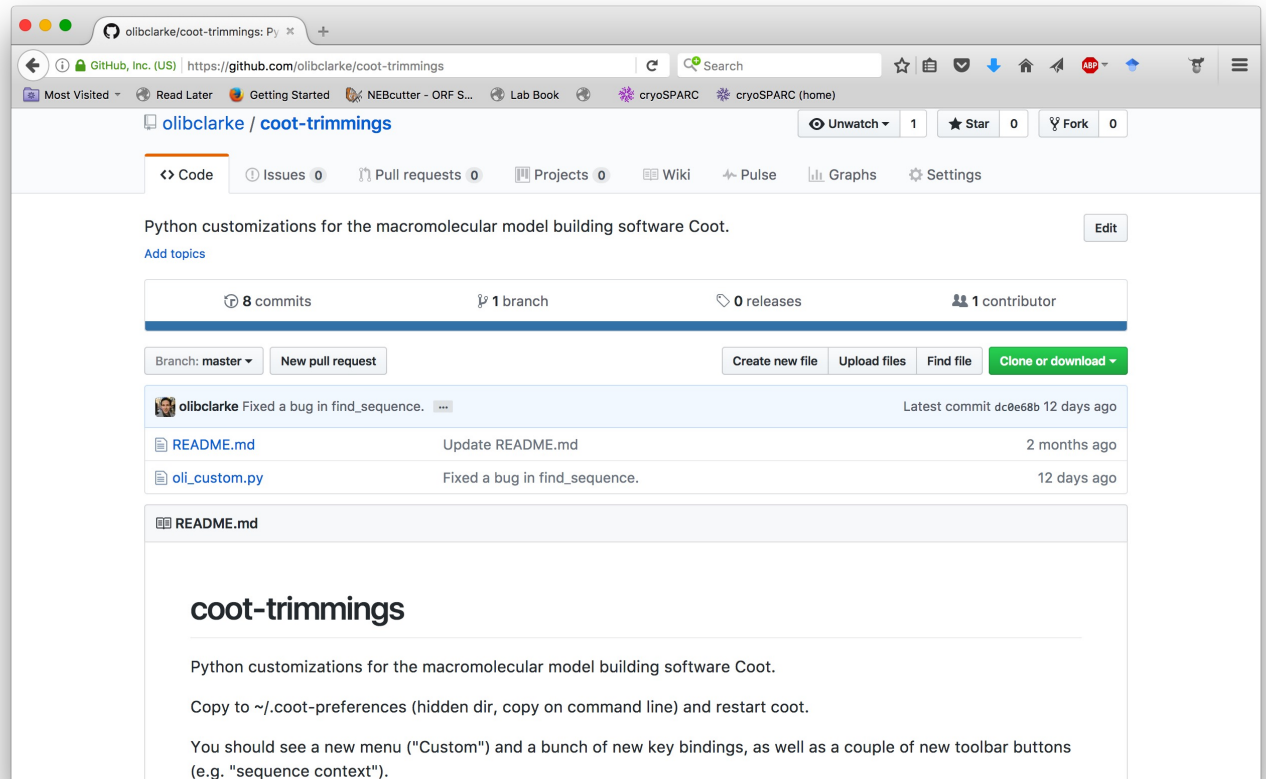- On-the-fly sharpening, resampling and low pass filtering

```python
def mutate_by_entered_code():
  def mutate_single_letter(X):
    entry=str(X).upper()
    mol_id=active_residue()[0]
    ch_id=active_residue()[1]
    resno=active_residue()[2]
    ins_code=active_residue()[3]
    resname=residue_name(mol_id,ch_id,resno,ins_code)
    map_id=imol_refinement_map()
    aa_dic={'A':'ALA','R':'ARG','N':'ASN','D':'ASP','C':'CYS','E':'GLU','Q':'GLN','G':'GLY','H':'HIS','I':'ILE','L':'LEU','K':'LYS
    nt_list=['A','C','T','G','U']
    if (resname in aa_dic.values()) and (aa_dic.get(entry,0)!=0):
      mutate(mol_id,ch_id,resno,ins_code,aa_dic.get(entry,0))
    elif (resname in nt_list) and (entry in nt_list):
      mutate_base(mol_id,ch_id,resno,ins_code,entry)
    else:
      info_dialog("Invalid target residue! Must be protein or nucleic acid, and entered code must be single letter.")
  generic_single_entry("New residue? (single letter code)","A","Mutate by single-letter code",mutate_single_letter)
```

```python
#mutate active residue to entered residue code (upper or lower case single-letter)
add_key_binding("Mutate by single letter code","M",
lambda: mutate_by_entered_code())
```

**Many pre-packaged functions available in COOT API. Mostly documented in online manual. Very easy to write your own! Useful e.g. for scripting domain-wise rigid body refinement.**
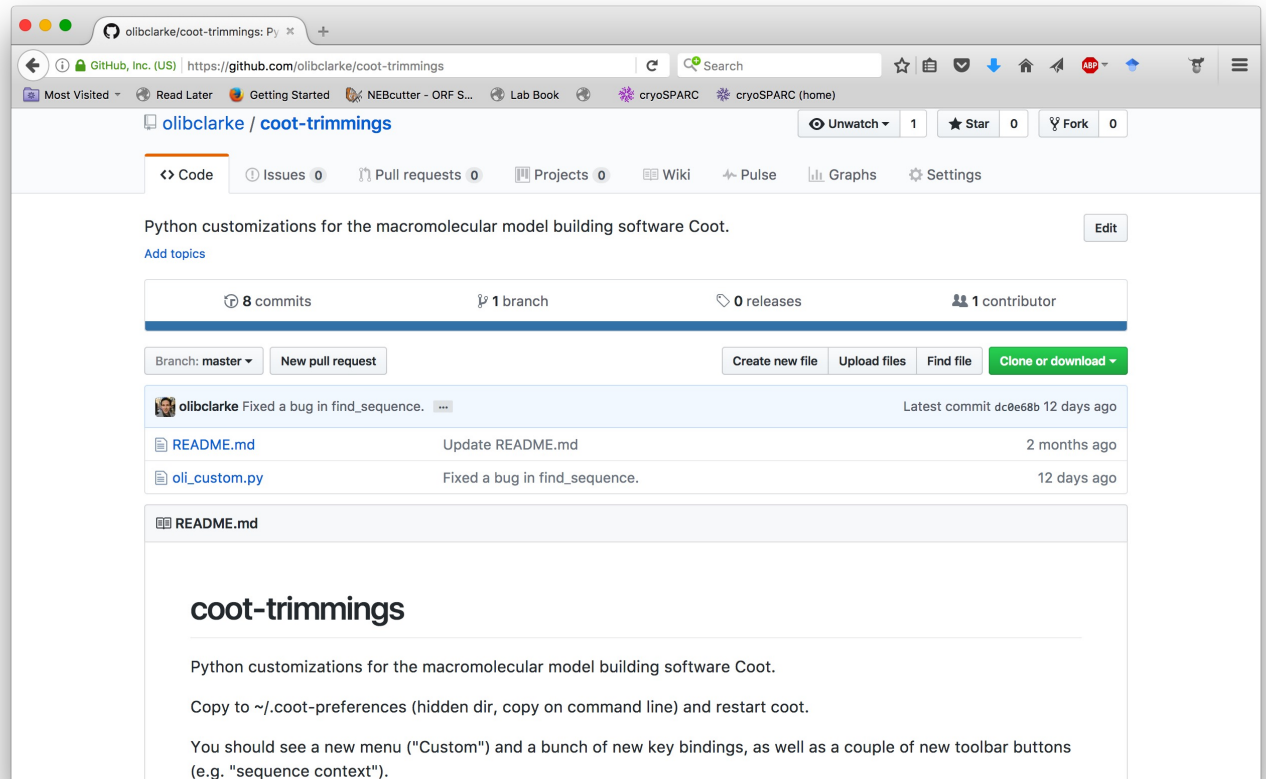
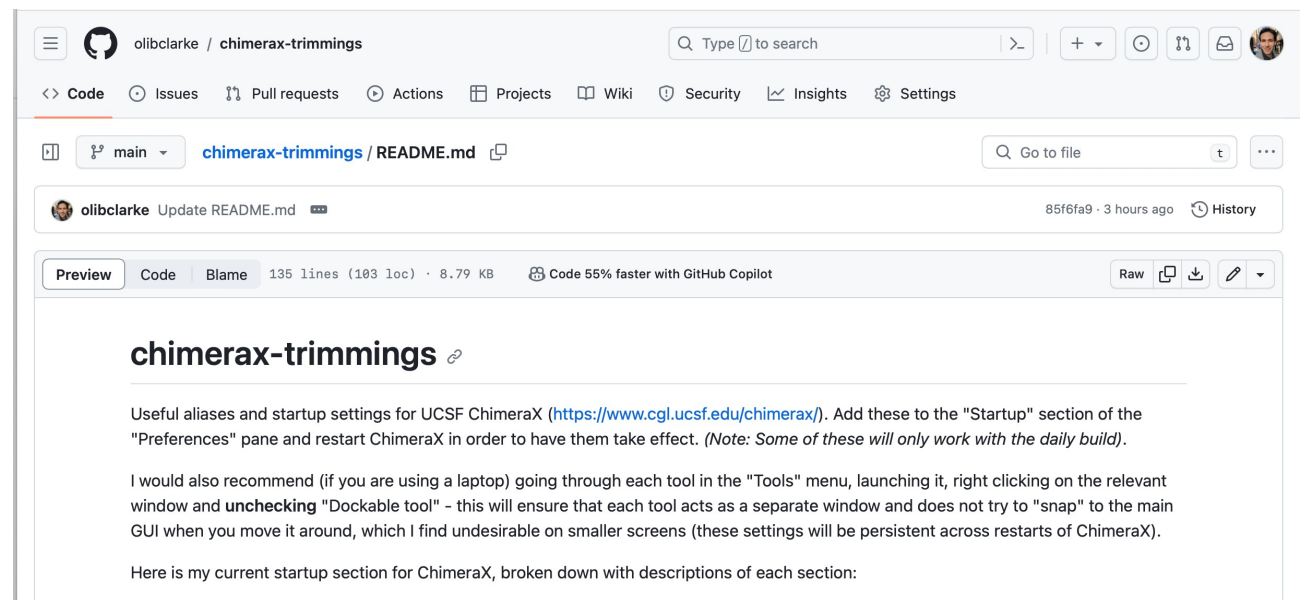**COOT – Crystallographic Object Oriented Toolkit**

- Simple, intuitive interface for building and manipulating atomic models in density maps.

- Low computational requirements

- **Extensive API – easy to script or modify (using simple Python code)**

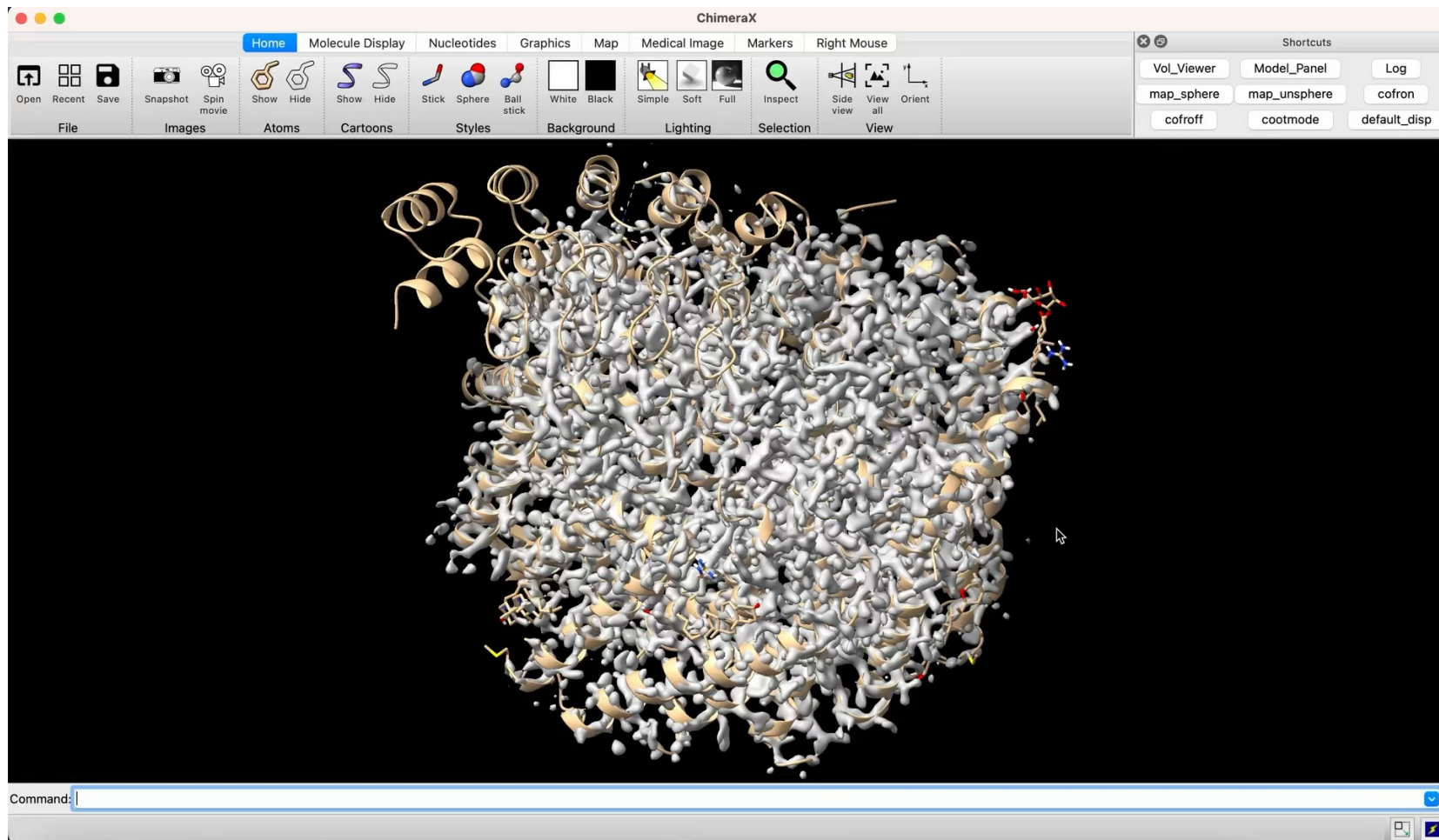- On-the-fly sharpening, resampling and low pass filtering



**Many pre-packaged functions available in COOT API. Mostly documented in online manual.**

**COOT – Crystallographic Object Oriented Toolkit**

- Simple, intuitive interface for building and manipulating atomic models in density maps.

- Low computational requirements

- **Extensive API – easy to script or modify (using simple Python code)**

- On-the-fly sharpening, resampling and low pass filtering



**Lots of key bindings, and easy to define custom keys. Learn them. They make everything much faster.**

**COOT – Crystallographic Object Oriented Toolkit**

- Simple, intuitive interface for building and manipulating atomic models in density maps.

- Low computational requirements

- **Extensive API – easy to script or modify (using simple Python code)**

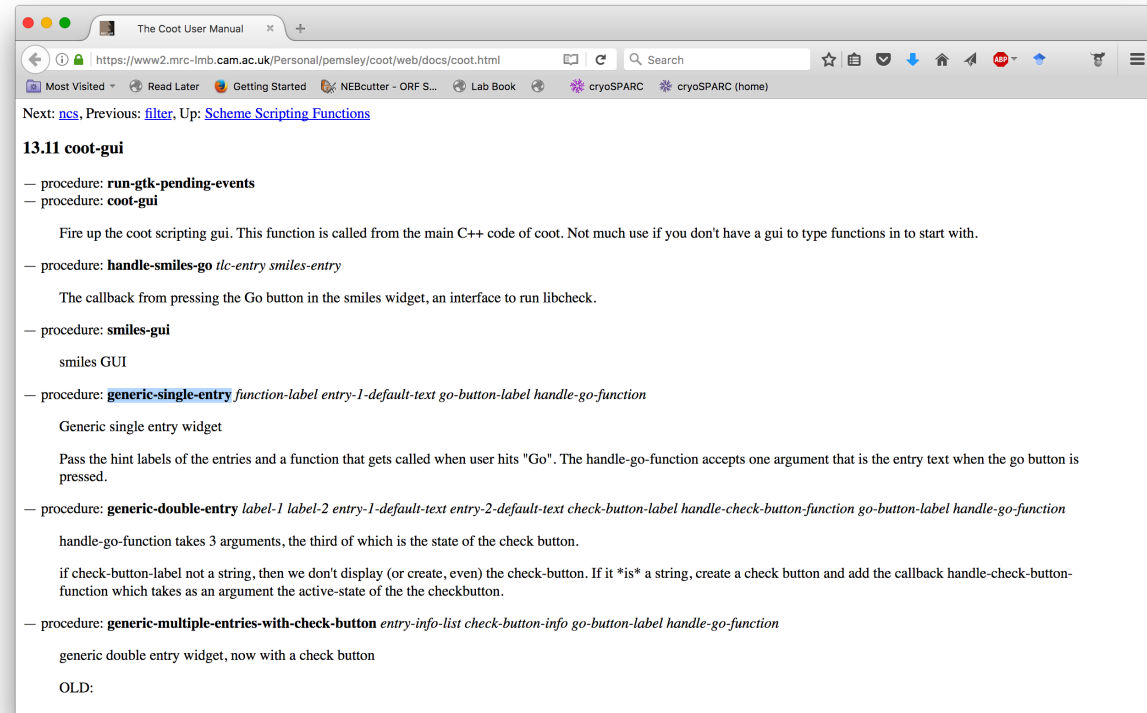- On-the-fly sharpening, resampling and low pass filtering



**Lots of key bindings, and easy to define custom keys. Learn them. They make everything much faster.**

## COOT – Crystallographic Object Oriented Toolkit

- Semi-automated helix placement

- Place cursor at the center of the helix and trigger "Place helix here" (I suggest via a key binding - "h" with *coot-trimmings*)

- Coot will attempt to automatically determine the length and direction of the helix.

- Trim/extend, adjust weights, then refine using real-space refine zone. Drag into density to adjust fit.

**COOT – Crystallographic Object Oriented Toolkit**

- Semi-automated helix placement

- **Place cursor at the center of the helix and trigger "Place helix here" (I suggest via a key binding - "h" with `coot-trimmings`)**

- Coot will attempt to automatically determine the length and direction of the helix.

- Trim/extend, adjust weights, then refine using real-space refine zone. Drag into density to adjust fit.

**COOT – Crystallographic Object Oriented Toolkit**

- Semi-automated helix placement

- Place cursor at the center of the helix and trigger "Place helix here" (I suggest via a key binding - "h" with `coot-trimmings`)

- Coot will attempt to automatically determine the length and direction of the helix.

- **Trim/extend, adjust weights, then refine using real-space refine zone. Drag into density to adjust fit.**

**COOT – Crystallographic Object Oriented Toolkit**

- Semi-automated helix placement

- Place cursor at the center of the helix and trigger "Place helix here" (I suggest via a key binding - "h" with `coot-trimmings`)

- Coot will attempt to automatically determine the length and direction of the helix.

- **Trim/extend, adjust weights, then refine using real-space refine zone. Drag into density to adjust fit.**

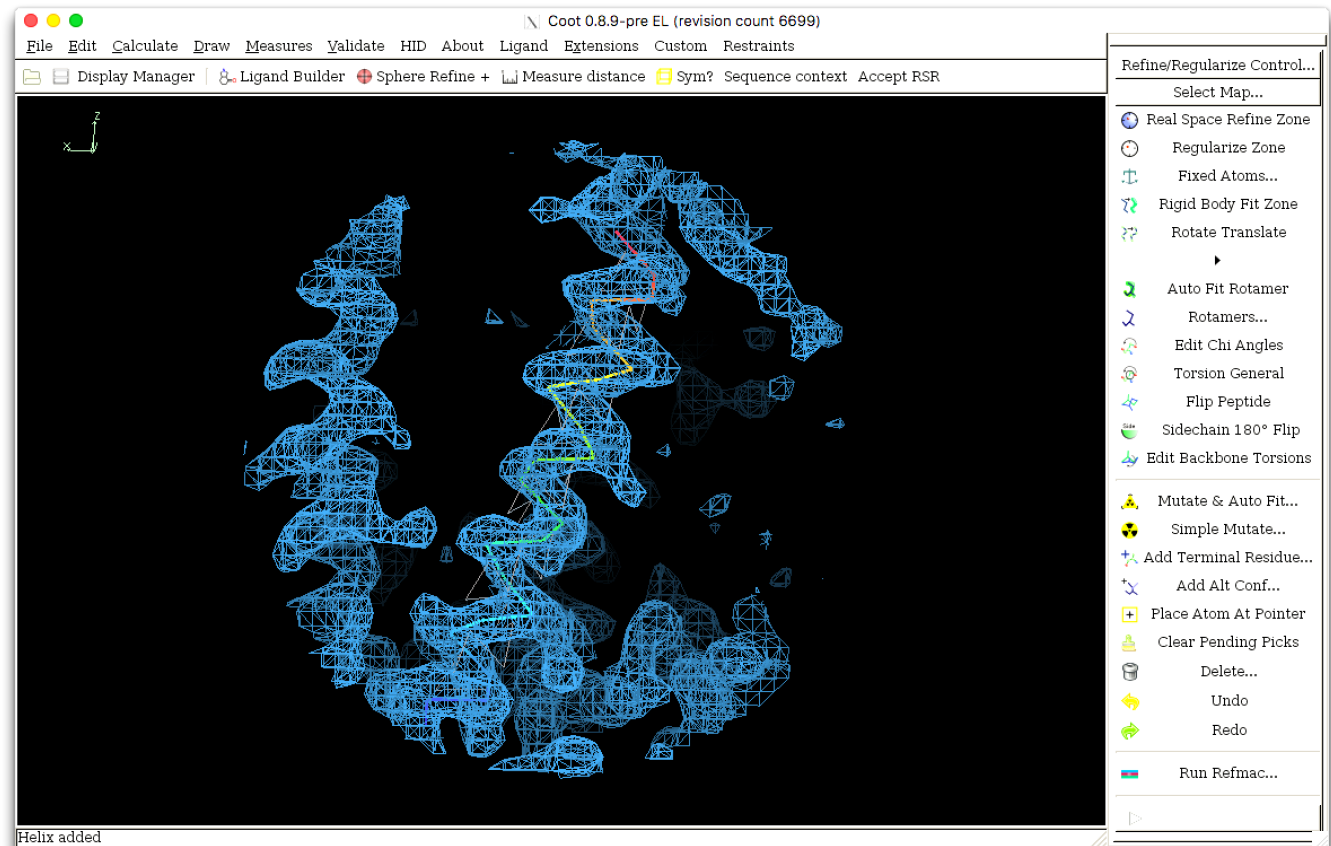## COOT – Crystallographic Object Oriented Toolkit

- Semi-automated helix placement

- Place cursor at the center of the helix and trigger "Place helix here" (I suggest via a key binding - "h" with `coot-trimmings`)

- Coot will attempt to automatically determine the length and direction of the helix.

- **Trim/extend, adjust weights, then refine using real-space refine zone. Drag into density to adjust fit.**
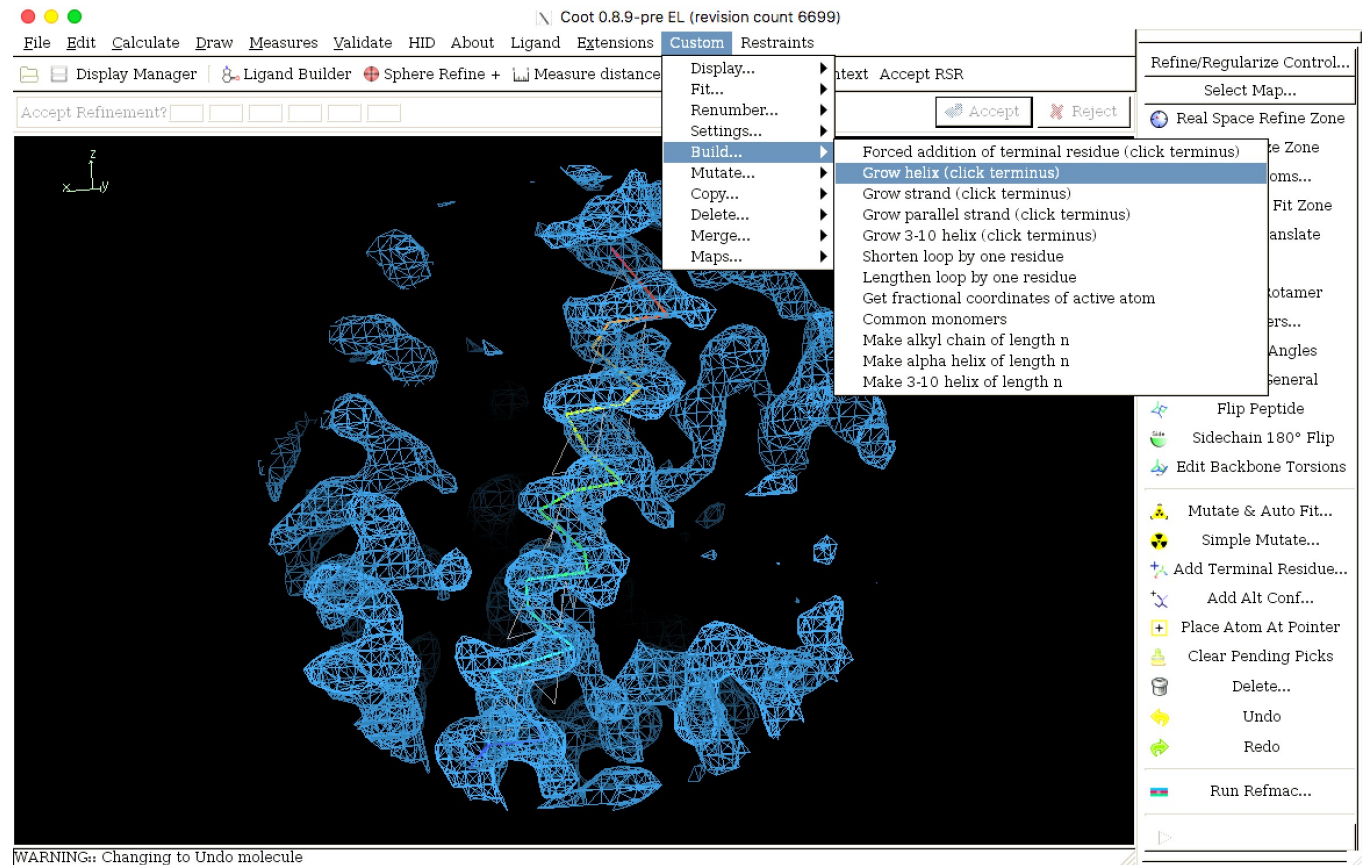
**COOT – Crystallographic Object Oriented Toolkit**

- Semi-automated helix placement

- Place cursor at the center of the helix and trigger "Place helix here" (I suggest via a key binding - "h" with *coot-trimmings*)

- Coot will attempt to automatically determine the length and direction of the helix.

- **Trim/extend, adjust weights, then refine using real-space refine zone. Drag into density to adjust fit.**
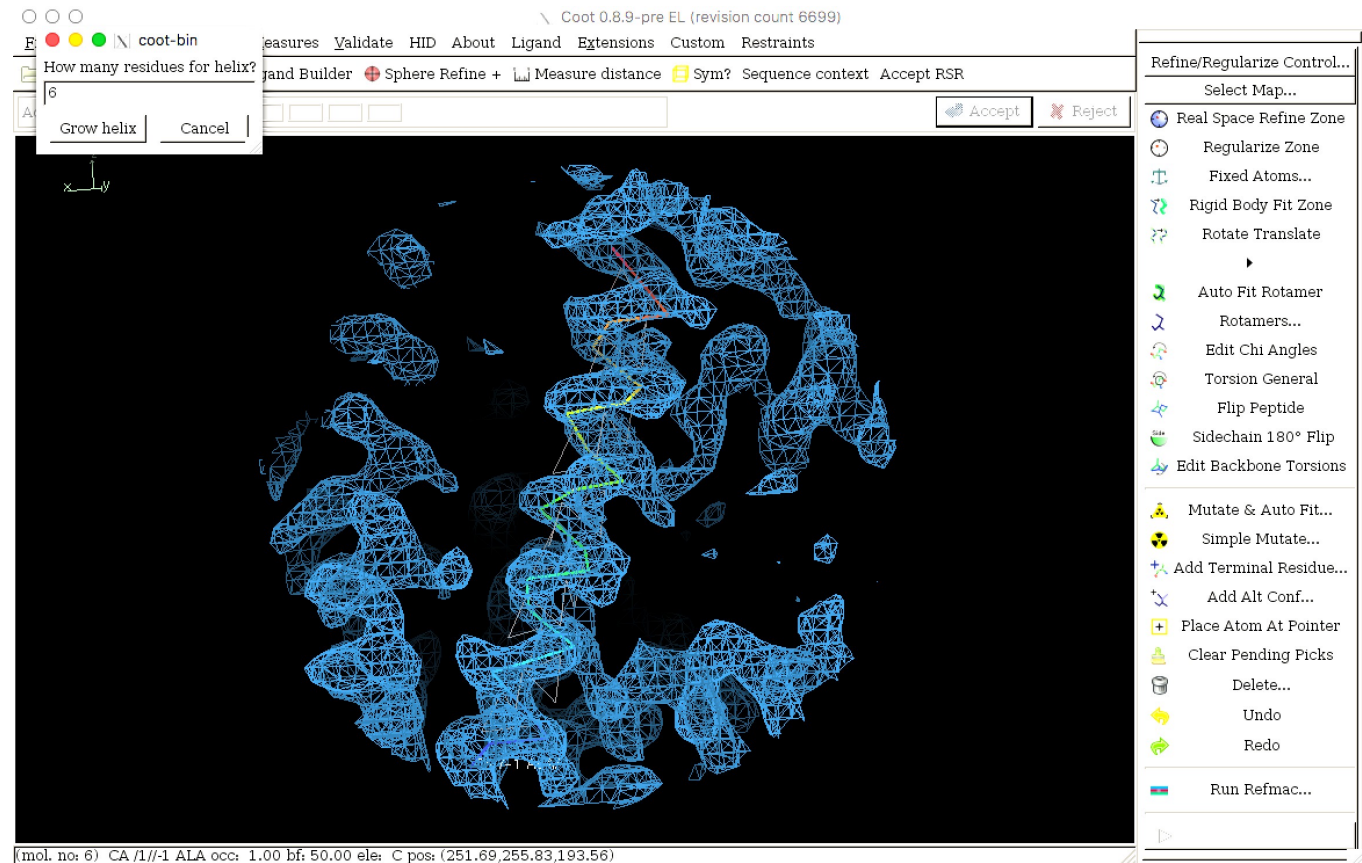
**COOT – Crystallographic Object Oriented Toolkit**

- Semi-automated helix placement

- Place cursor at the center of the helix and trigger "Place helix here" (I suggest via a key binding - "h" with `coot-trimmings`)

- Coot will attempt to automatically determine the length and direction of the helix.

- **Trim/extend, adjust weights, then refine using real-space refine zone. Drag into density to adjust fit.**
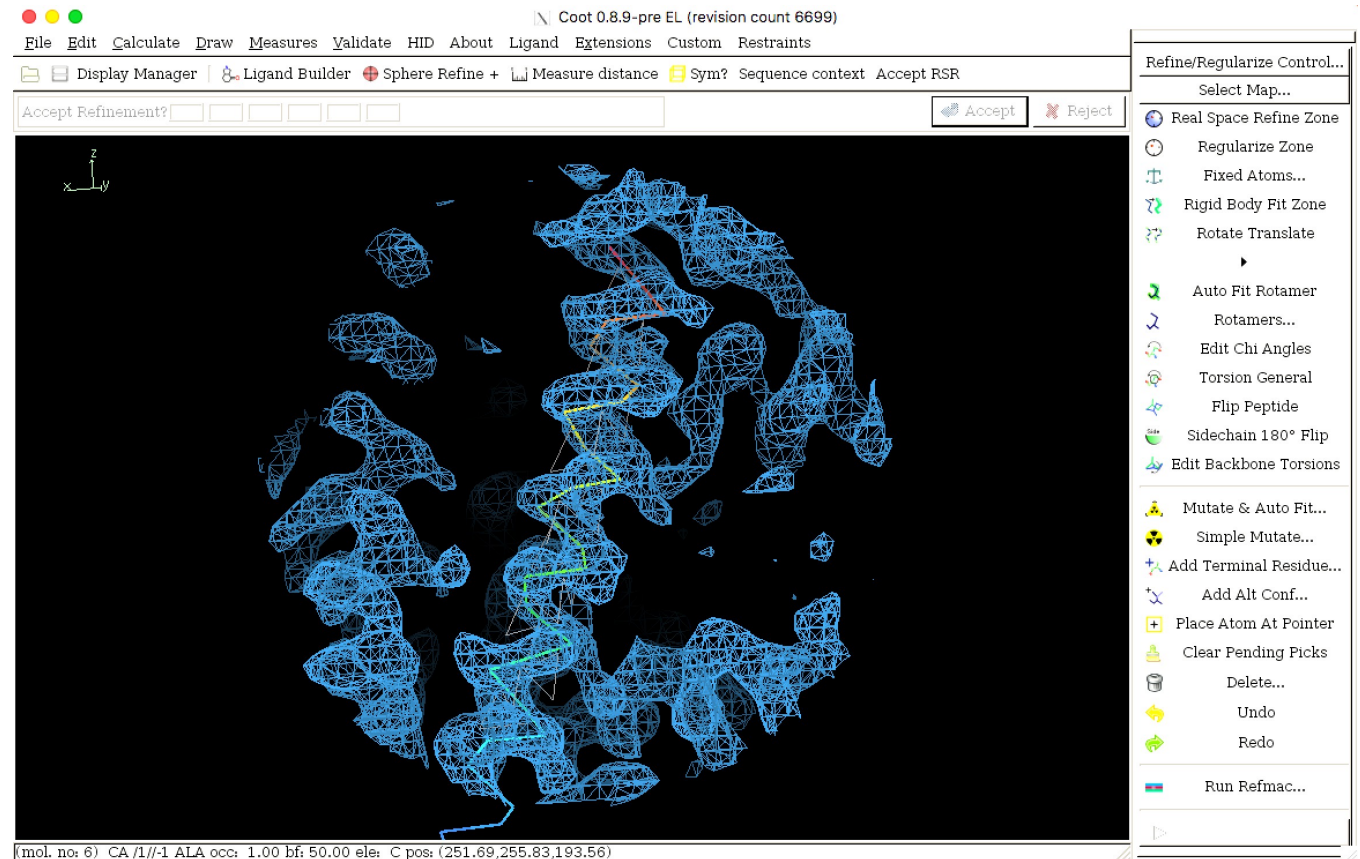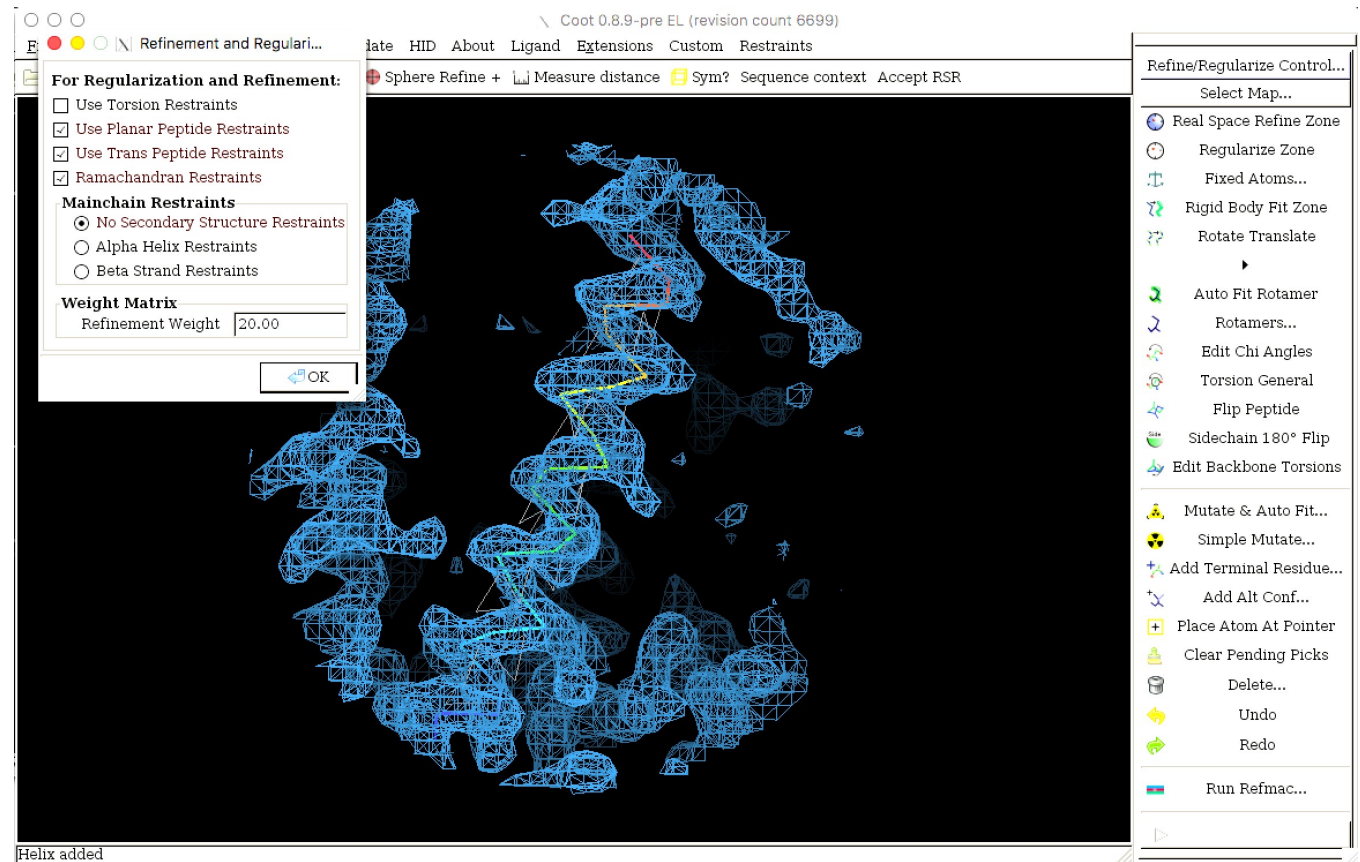
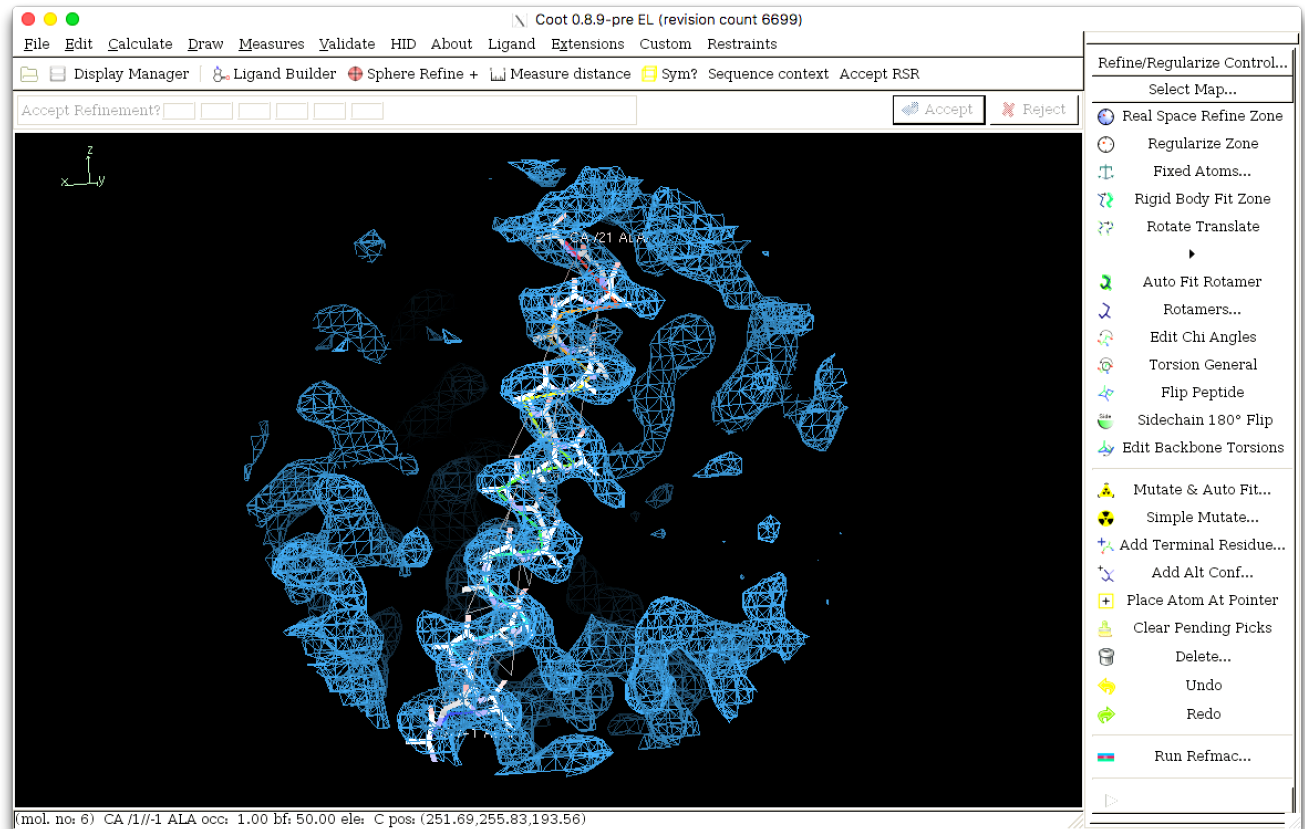**COOT – Crystallographic Object Oriented Toolkit**

- Sequence assignment.

- **Adjust numbering to match expected position in sequence.**

- Mutate to match sequence

- Fill sidechains manually.

- Adjust sequence register to optimize local fit to sidechain densities.

**COOT – Crystallographic Object Oriented Toolkit**

- Sequence assignment.

- **Adjust numbering to match expected position in sequence.**

- Mutate to match sequence

- Fill sidechains manually.

- Adjust sequence register to optimize local fit to sidechain densities.

**COOT – Crystallographic Object Oriented Toolkit**

- Sequence assignment.

- Adjust numbering to match expected position in sequence.

- **Mutate to match sequence**

- Fill sidechains manually.

- Adjust sequence register to optimize local fit to sidechain densities.

**COOT – Crystallographic Object Oriented Toolkit**

- Sequence assignment.

- Adjust numbering to match expected position in sequence.

- **Mutate to match sequence**

- Fill sidechains manually.

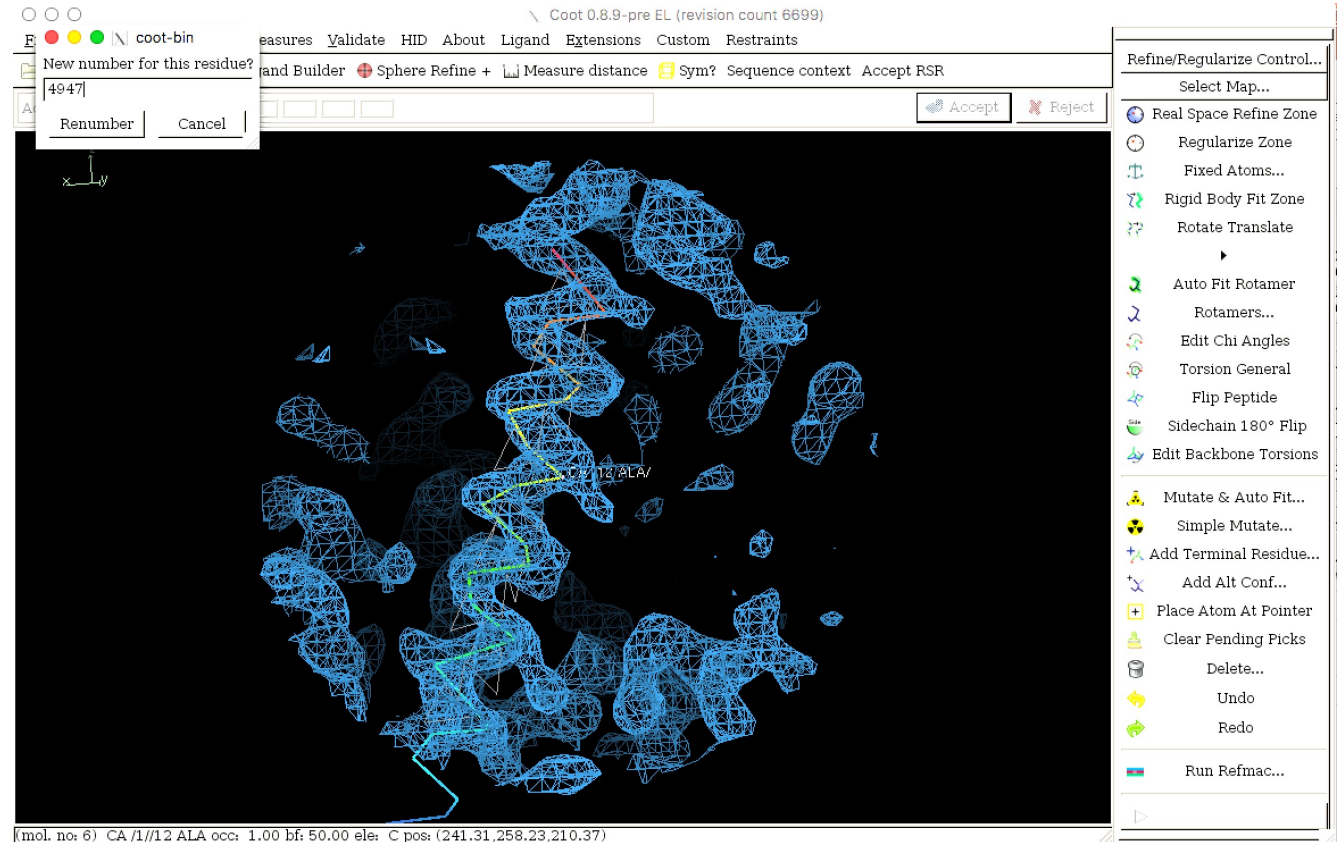- Adjust sequence register to optimize local fit to sidechain densities.

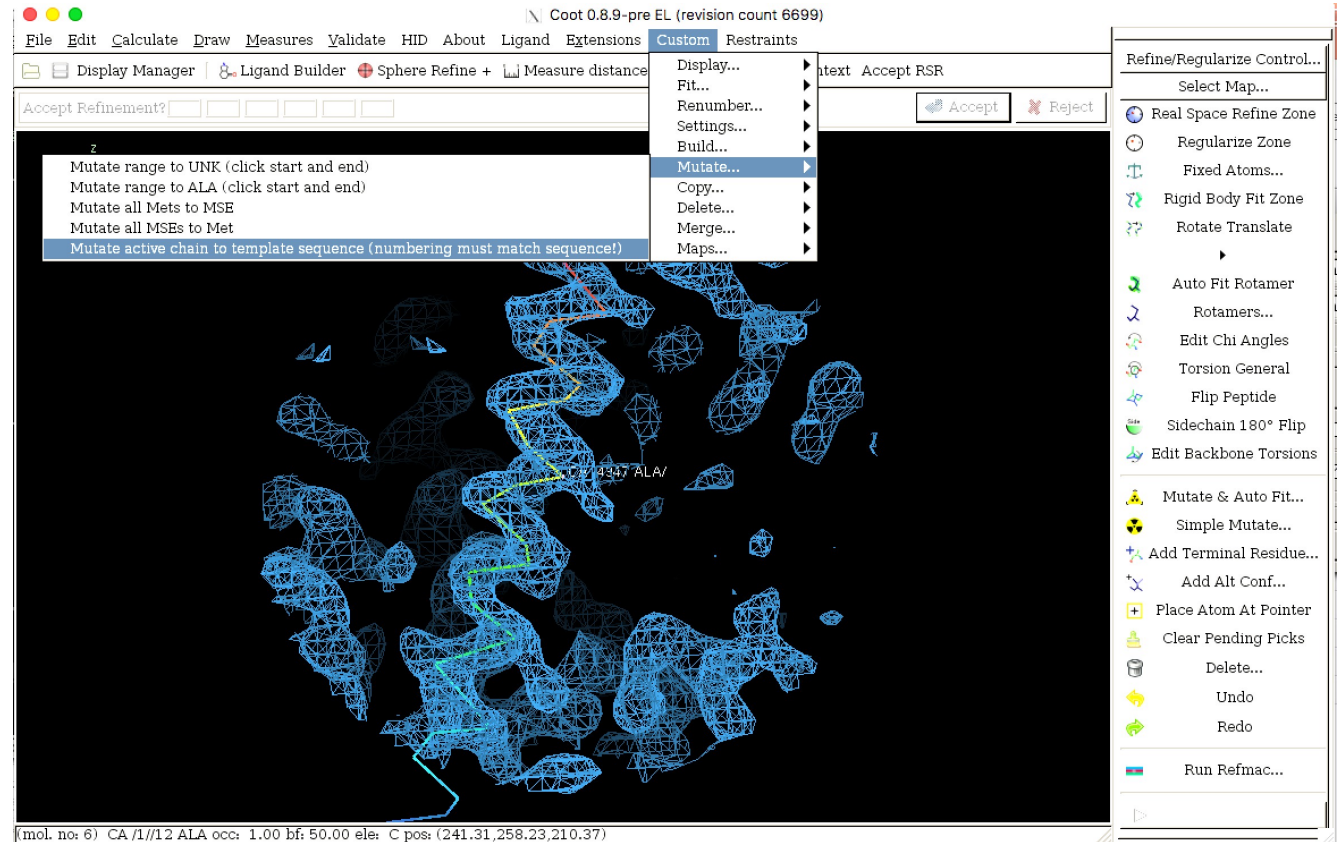## COOT – Crystallographic Object Oriented Toolkit

- Sequence assignment.

- Adjust numbering to match expected position in sequence.

- **Mutate to match sequence**

- Fill sidechains manually.

- Adjust sequence register to optimize local fit to sidechain densities.
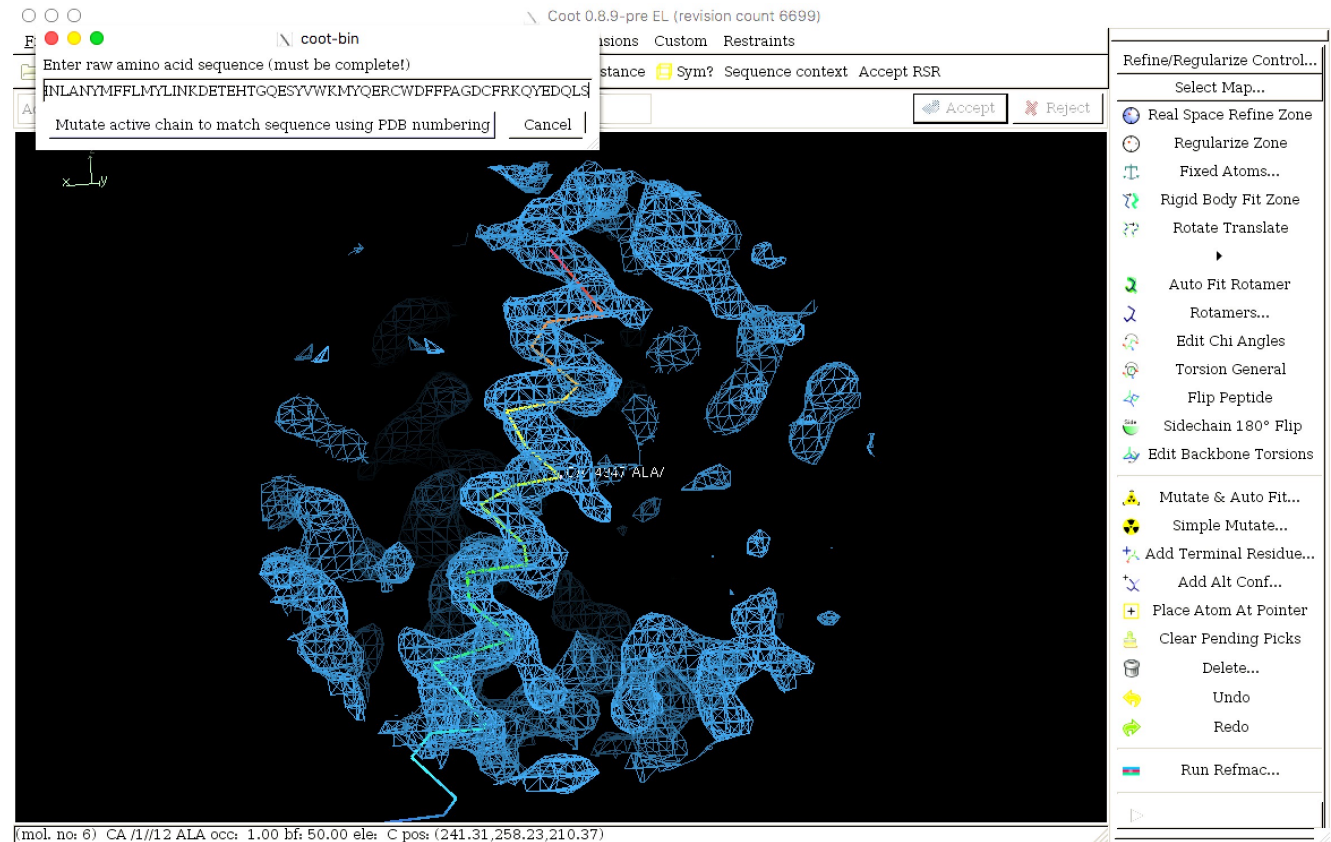


**Use 'Add Terminal residue' to extend chain.**

**COOT 1.0 coming soon! Improved visualization, refinement, modern GUI**

**COOT 1.0 coming soon! Improved visualization, refinement, modern GUI**



Contour using one map,
Colour using another

**COOT 1.0 coming soon! Improved visualization, refinement, modern GUI**

**COOT 1.0 coming soon! Improved visualization, refinement, modern GUI**

**COOT 1.0 coming soon! Improved visualization, refinement, modern GUI**

**COOT 1.0 coming soon! Improved visualization, refinement, modern GUI**

**COOT 1.0 coming soon! Improved visualization, refinement, modern GUI**

## ISOLDE

- Interactive molecular dynamics flexible fitting, implemented as plugin for ChimeraX

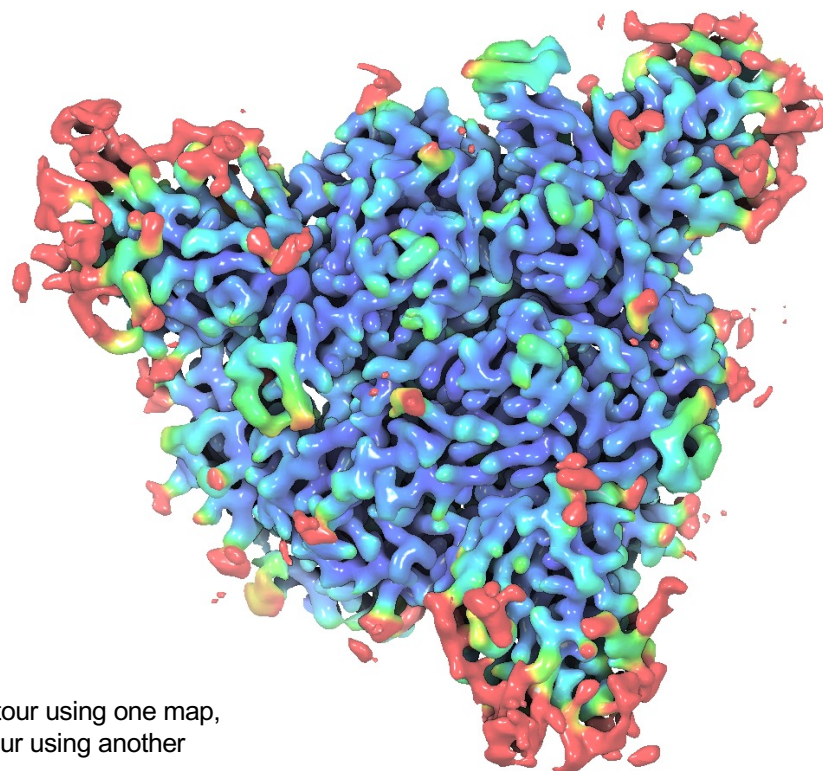- Useful during "polishing" stage of generating a final model, identifying and fixing otherwise difficult to correct errors in geometry, non-bonded contacts. Physically realistic simulation guided by map, user input.

- Complementary to COOT – COOT better for de novo building and assembly, ligand placement, ISOLDE very useful for final round of real space fitting.

*(Croll, 2018, Acta. Cryst. D)*

# Types of errors in macromolecular models

- Identity (e.g. wrong domain)

- Directionality

- Topology/connectivity

- Register

- Rotamer

- Backbone torsion

- Ligand identification and placement

**Types of errors in macromolecular models**

- Identity (e.g. wrong domain)

- Directionality                                    Low resolution (<4.5 Å)

- **Topology/connectivity**

                                                    Medium resolution (3.5-4.5 Å)
- **Register**

- Rotamer

- Backbone torsion                    Medium/high resolution (2.5-4 Å)

- Ligand identification and placement

**Types of errors in macromolecular models**

- Identity (e.g. wrong domain)

- Directionality

- Topology/connectivity

- **Register**

- Rotamer

- Backbone torsion

- Ligand identification and placement

**Types of errors in macromolecular models**

- Identity (e.g. wrong domain)

- Directionality

- Topology/connectivity

- **Register**

- Rotamer

- Backbone torsion

- Ligand identification and placement

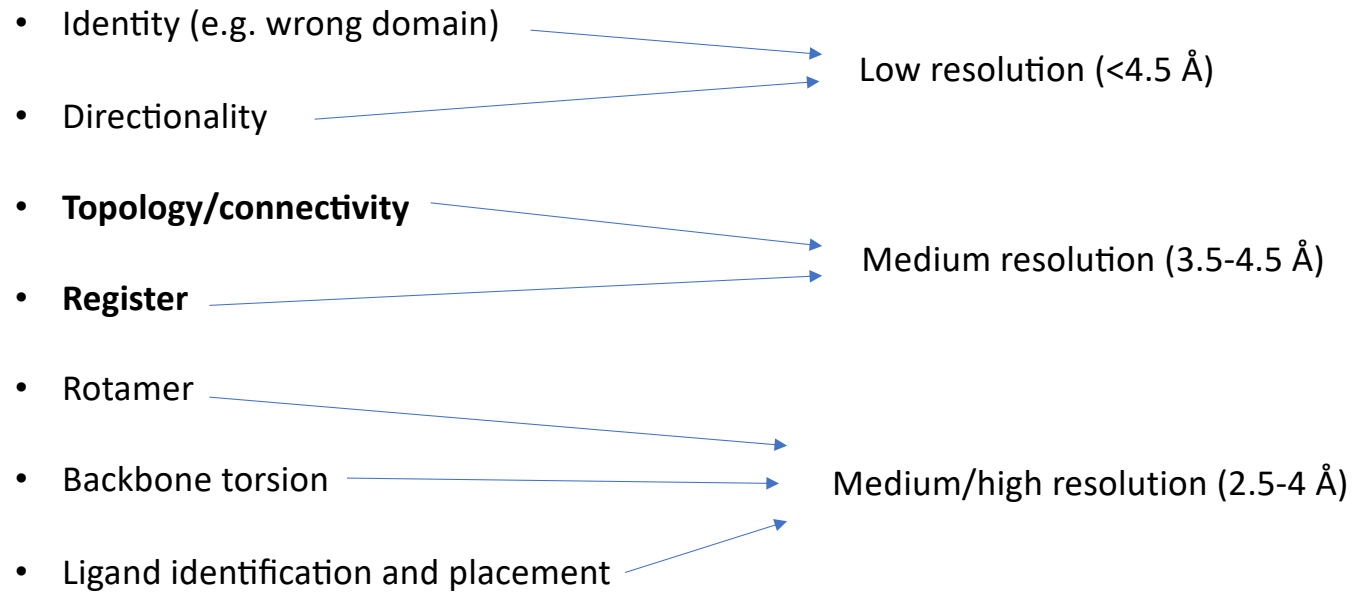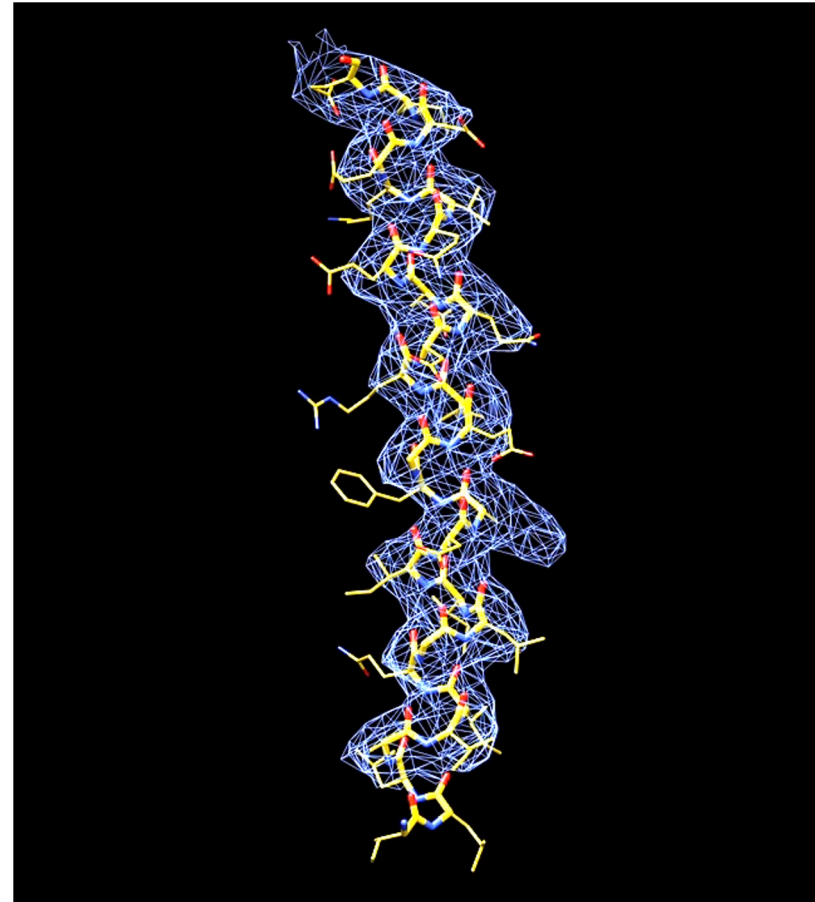**Types of errors in macromolecular models**

- Identity (e.g. wrong domain)

- Directionality

                                        Low resolution (<4.5 Å)

- Topology/connectivity

- Register

                                        Medium resolution (3.5-4.5 Å)

- Rotamer

- Backbone torsion

                                        Medium/high resolution (2.5-4 Å)

- **Ligand identification and placement**

**Cautionary note – don't believe everything you read in the (Nature/Science/Cell) paper…**



Orthosteric site

Gαi
scFv16
Gγ
Gβ

Gαiq
scFv16
Gγ
Gβ

**Cautionary note – don't believe everything you read in the (Nature/Science/Cell) paper…**



Ligand invisible at any reasonable contour; protein density anisotropic and broken.
So how was nice density figure for ligand made? By carving the map around the ligand. Please don't do this!

**Cautionary note – don't believe everything you read in the (Nature/Science/Cell) paper…**



**Ligand invisible at any reasonable contour; protein density anisotropic and broken.**
**So how was nice density figure for ligand made? By carving the map around the ligand. Please don't do this!**
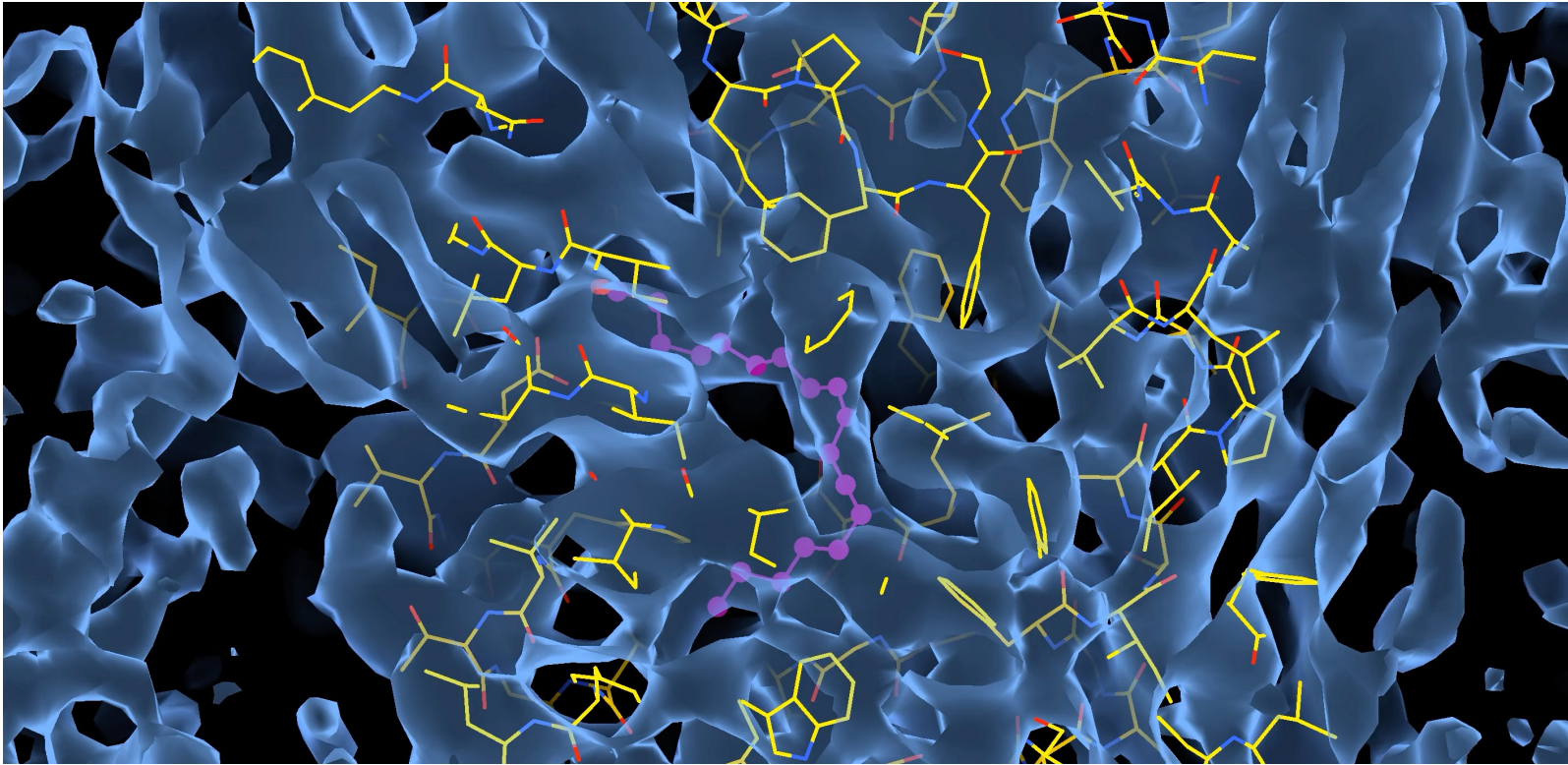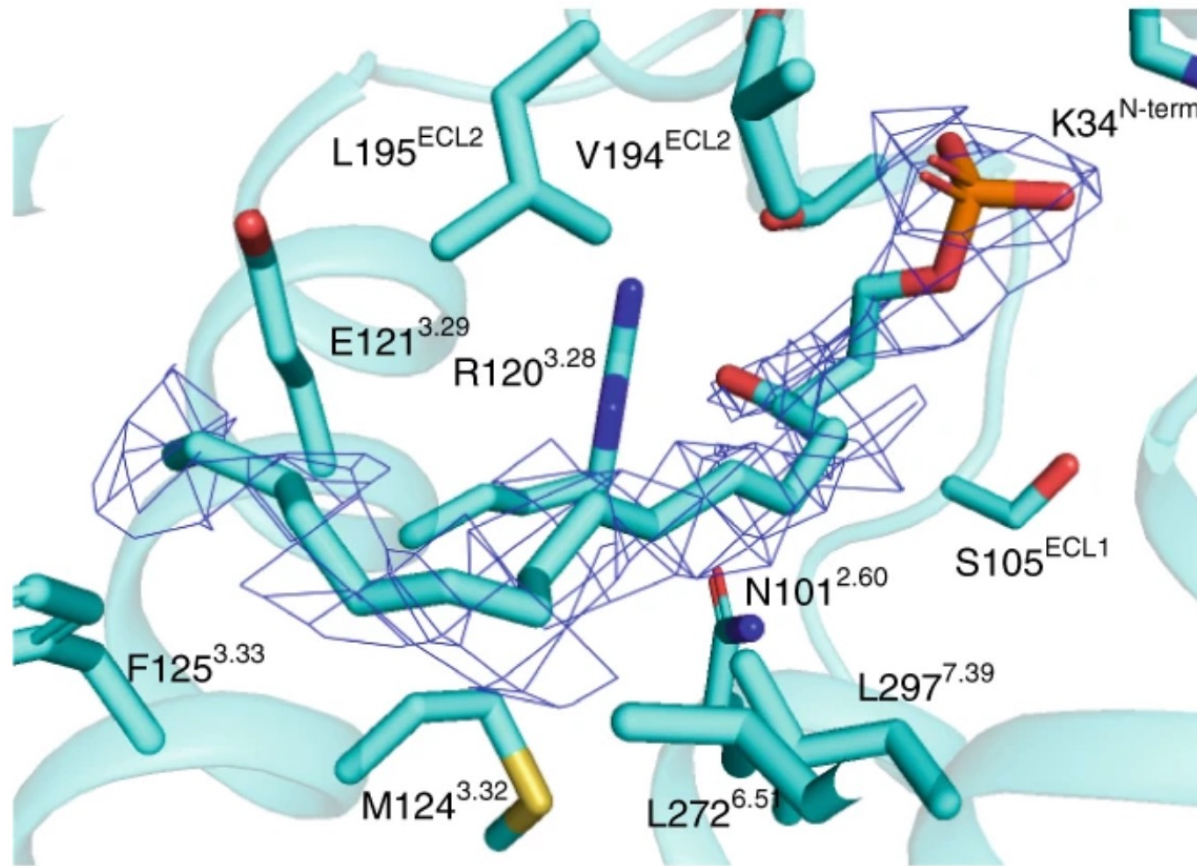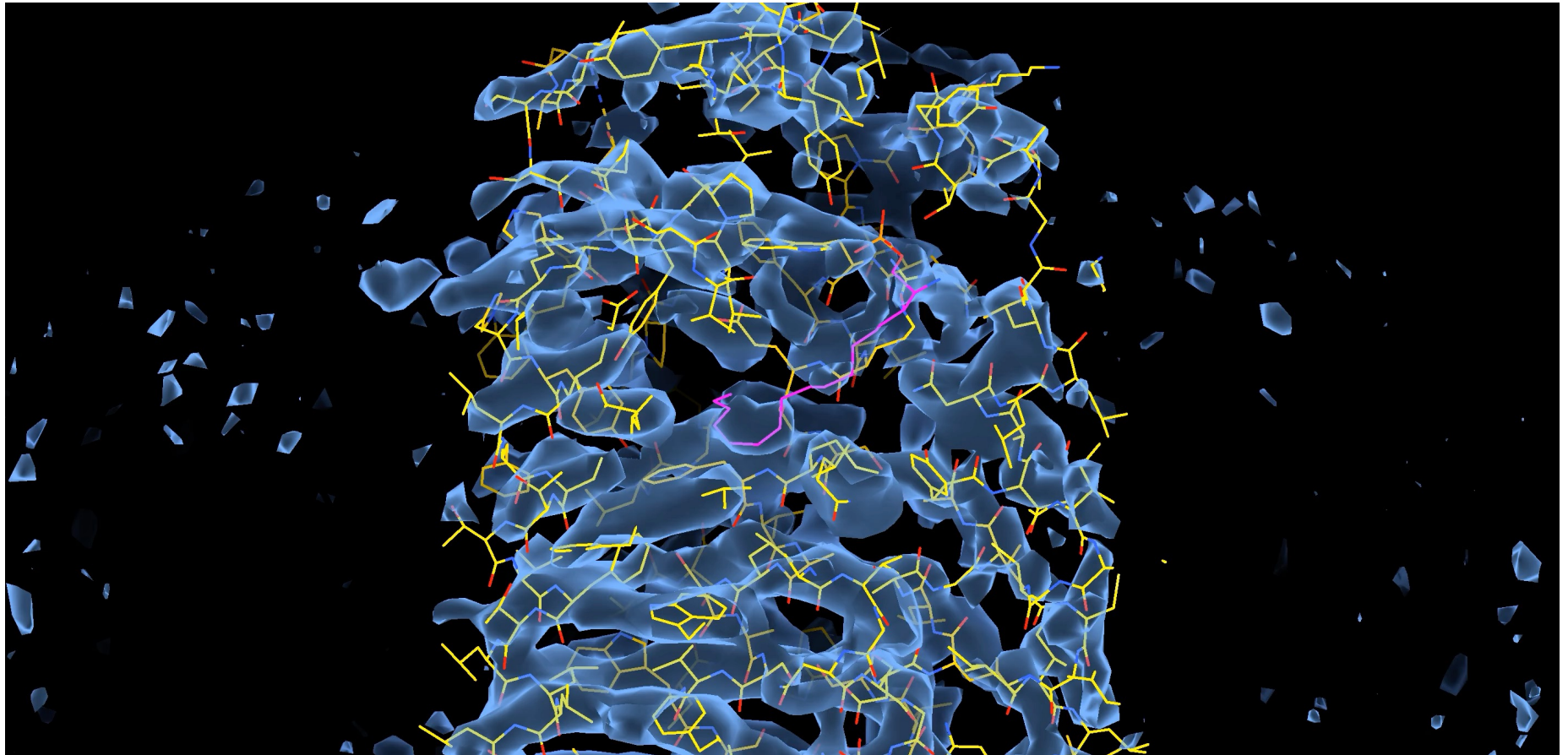
**Cautionary note – don't believe everything you read in the (Nature/Science/Cell) paper…**

**Cautionary note – don't believe everything you read in the (Nature/Science/Cell) paper...**

**Strategy for identifying and correcting errors.**

- Analyse as you go – "sanity checks" on chemistry, nonbonded interactions, surface composition. Use Molprobity for clashes, Chimera(X) or pymol to check e.g. for buried polars, exposed hydrophobics. Monitor agreement with secondary structure, disorder predictions.

- Use EM-ringer (or other local metrics – Q-score, Strudel score, MEDIC all useful) to identify errors in backbone and rotamer geometry.

- **Look at everything! Manually check and recheck the fit of every residue. Tedious but necessary.**

- Sometimes, you just can't tell the right answer. Don't be afraid to specify sequence ambiguity (use UNKs).

- Half-map FSCs are useful to analyze/identify overfitting during refinement, but they tell you little about the local quality or correctness of the model.

**Finally…**

**"ALL MODELS ARE WRONG, BUT SOME ARE USEFUL"** – *George P. Box*

\* It should be remembered that just as the Declaration of Independence promises the <u>pursuit</u> of happiness rather than happiness itself, so the iterative scientific model building process offers only the pursuit of the perfect model. For even when we feel we have carried the model building process to a conclusion some new initiative may make further improvement possible. Fortunately to be useful a model does not have to be perfect.

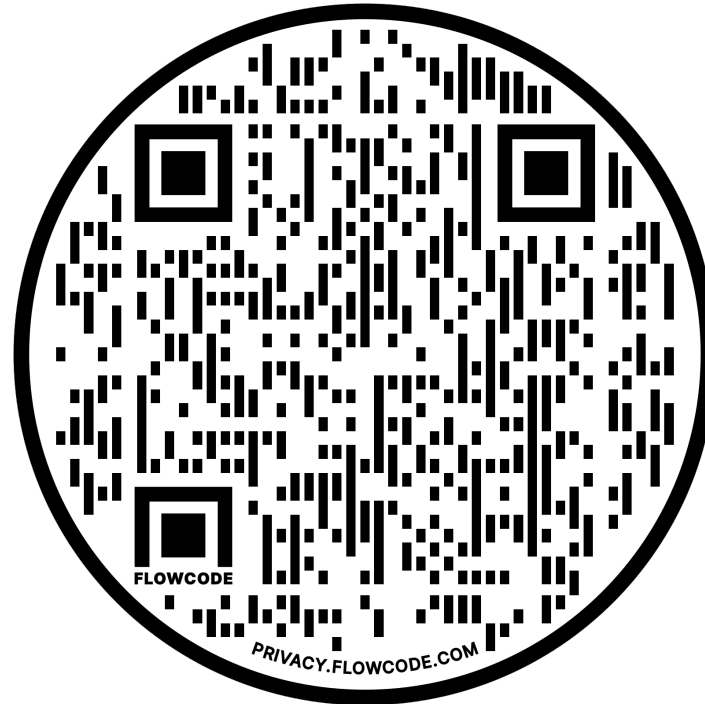*George P. Box, "Robustness in the Strategy of Scientific Model Building", 1979*

**Thank you for listening!**

Columbia University
Medical Center

**Model building tutorial**



Tutorial PDF: https://bit.ly/2XPsiox

Data: https://bit.ly/3ASQ41I

AlphaFold add on: https://bit.ly/3KTo6qX