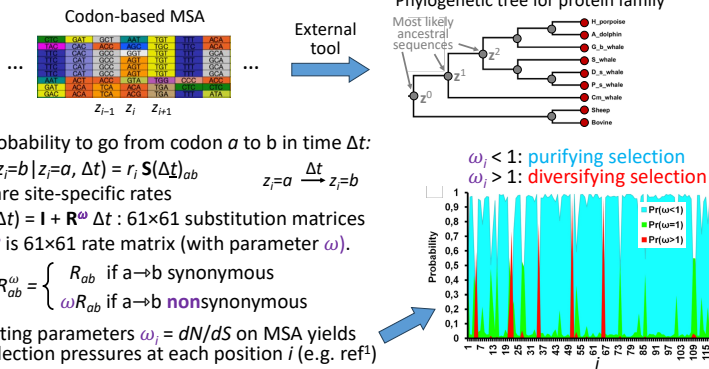


C2

## Novel methods for analysis of selection using mutation-selection

Johannes  
Söding

### State of the art: codon substitution models



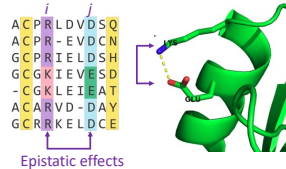
### Objective: More realistic model to analyse selection on proteins<sup>2,3</sup>

- Model *nearly neutral theory of evolution*
- *Epistatic fitness function*
- *Heterogeneous codon frequencies*
- *Rate heterogeneity* as emergent property<sup>2</sup>
- *Heterotachy* as emergent property<sup>2</sup>
- *GC bias* (in nucleotide mutation matrix  $\mu_{ab}$ )
- Estimate selection pressure ( $\sim dN/dS$ ) for each protein position and each tree node *relative to rest of phylogenetic tree*

### PhD 1 - Inference method for epistatic MutSel model of molecular evolution

Goal: Develop 'realistic' epistatic MutSel models, train and use to sample ancestral sequences

Proteins need to maintain stably folding structure. This requires many residue-residue interactions. So, fitness effect of mutations depends strongly on entire sequence  $\Rightarrow$  **strong epistasis**



MutSel:  $p(z' | z, \Delta t) = \mu_{ab} \Delta t p_{\text{fix}}(s_{z \rightarrow z'})$

Fixation probability:  
 $p_{\text{fix}}(s_{z \rightarrow z'}) = \frac{1 - \exp(-2 s_{z \rightarrow z'})}{1 - \exp(-2N s_{z \rightarrow z'})}$   $N = \text{effective population size}$

Selection coefficient:  $s_{z \rightarrow z'} = f_{v,w}(z') - f_{v,w}(z)$

Fitness function:  $f_{v,w}(z) \propto \left( \exp \left[ \sum_i v_{i0} [z_i=a] + \sum_{i < j} w_{i0,jb} [z_i=a, z_j=b] \right] \right)$  Epistatic effects

Learn  $20L + 20^2L^2$  parameters  $v$  and  $w$  from large **protein language model** (e.g. MSA transformer), fix  $w$ , fit only  $v$  on MSA. (using *stochastic variational inference* and *Gibbs sampling*)

**Sampling** ancestral sequences allows for computing uncertainties.

**GC bias**: set  $\mu_{ab}^s$  such that equilibrium GC frequency is equal to mean in codon sequences of **species  $s$** . Estimate for ancestors.

### PhD 2 - Methods for analysis of selection based on MutSel model with epistatic fitness

Goal: New Bayesian methods for position- and time-resolved analysis of selection on proteins

Current tools for MutSel-based selection analysis not much used

- Not well calibrated: they always yield  $\omega \leq 1$  (no pos. sel.!)<sup>4</sup>
- No time-resolved analysis over tree nodes
- No epistasis

Calibration is a conceptual issue with common MutSel tools. e.g.

$$\omega_i = dN_i/dS_i = N/N_{\text{mut}} = \frac{\sum_n \sum_a \sum_b \text{non-syn. } \mu_{ab} \Delta t p_{\text{fix}}(s_{a \rightarrow b}) p(z_i^n=a)}{\sum_n \sum_a \sum_b \text{non-syn. } \mu_{ab} \Delta t p(z_i^n=a)} \leq 1$$

Here

- Include estimates of uncertainties by posterior sampling over ancestral sequences
- Assess selection pressure for each node  $n$  and position  $i$
- Calibrate  $\omega$  values relative to (constant) fitness model over entire tree:

$$\omega_{ni} = dN_{ni}/dS_{ni} = \frac{\mathbb{E}_z \left[ \sum_a \sum_b \text{non-syn. } p(z_i^m=b | z_i^n=a, \Delta t) \right]}{\mathbb{E}_z \left[ \sum_a \sum_b \text{syn. } p(z_i^m=b | z_i^n=a, \Delta t) \right]} = \frac{\mathbb{E}_z \left[ \sum_a \sum_b \text{non-syn. } p(z_i=b | z_i^n=a, \Delta t) \right]}{\mathbb{E}_z \left[ \sum_a \sum_b \text{syn. } p(z_i=b | z_i^n=a, \Delta t) \right]}$$

Expectation values over many samples of ancestral sequences

Realizations of  $z_i^m$  given leaf nodes

Realizations of  $z_i^m$  ignoring leaf nodes

- Open source software

### Collaborations with A1, A2, A4, B3, B4, B5

- Develop visualizations / analysis plots "on the job"
- Improve and extend evolutionary model

### Collaborations with all RTG projects

- Assistance with deep annotation and seq/structure analysis (Foldseek, Collabfold, HHpred, MMseqs,...)

### References

1. Dasmeh P et al. 2013. Positively Selected Sites in Cetacean Myoglobins Contribute to Protein Stability. PLOS Comput Biol 9:e1002929.
2. de la Paz J et al. 2020. Epistatic contributions promote the unification of incompatible models of neutral molecular evolution. PNAS 117: 5873-82.
3. Vorberg S, Seemayer S, Söding J. 2018. Synthetic protein alignments quantify noise in residue-residue contact prediction. PLOS CB 14 e1006526.
4. Spielman SJ, Wilke CO. 2015. The relationship between dN/dS and scaled selection coefficients. MBE 32:1097-1108.

RTG  
2984/1