

Structural and functional constraints on protein evolution

Claus O. Wilke

The University of Texas at Austin

... P N G E N R K - N I E K F T E K K N F F Y R L K K F I E K N F - S P N D F V I L L M G D I N V A I N D K D I G I S E N N ...
... P Q G E N R N - N K I K F N K K K N F Y K N L I K Y V Q K K I - F E N K N I I I L M G D M N I S P E D Q D I G I D P Y H ...
... P H G E S F Y - K T D K F E E K K F F Y Q K L Y L F L K R N C - E K N A H I L I M G D M N I S P T D L D I G L S V N S ...
... P N G E S K K - N L K K F E I K K I F Y K S L F L F L N K F Y - E K N M H I L I M G D M N I S P E D L D V G L S L T S ...
... P Q G K S I D -- H P D Y Q A K H R F F D R L L N L F Q R E F - S P H T P L L W V G D M N V A P T D I D V --- T S ...
... P N G N P A P -- G P K F D Y K L R W F E R L R L R A Q E L I - A T G A P V V I A G D Y N V M P T E L D V --- Y K ...
... P N G N P A P -- G P K F D Y K L R W F D R L I T H A Q G L L - A A G K P V L L T G D F N V M P T E L D V --- Y K ...
... P N G N P V E -- T P K Y P Y K L R W M D R L I R Y A E D R L - A L E E P L V L A G D Y N V L P T P D D V --- A N ...
... P N G N P P Q -- T E K Y P Y K L K W M D R L L A Y S K E R L - K S E E P F V L A G D F N V I P T P E D V --- Y N ...
... P N G N P P N -- T E K Y P Y K L K W M S R L R D Y A R E R L - K T E E P L I L A G D F N V I P A A A D V --- S N ...
... P N G N P V G -- S E K Y P Y K L S W M A R L R D Y A Q Q R L - K T E E P L I L A G D F N V I P Q A E D V --- H N ...
... P N G N P V P -- G P K Y D Y K L A W M E R L R A R A I E L L - K S E A P F V M A G D Y N I I P Q P M D A --- A K ...
... P N G N P A P -- G P K Y D Y K L A W M A R M H A R V E S L L - P L E E P L V F C G D Y N V I P Q A E D A --- A K ...
... P N G N P A P -- G P K Y D Y K L A W M E R L E A R A R E E L L - A E E M P A L M A G D Y N V I P Q A E D A --- A R ...
... P N G N P A P -- G P K Y D Y K L A W M E R L R A R A E A L L - K A E E P A L M A G D Y N V I P Q A E D A --- A K ...
... P N G N P A P -- G P K F D Y K L A W M Q R L E A R A K A L L - A D E M P F I L M A G D Y N I I P Q A E D A --- A K ...
... P N G N P V D -- T E K F S Y K L E W M D R L I A R A K E L L - L L E E P F V M M G D Y N I I P H E D D V --- H D ...
... P N G N P I D -- S D K F P Y K L S W M E R L R S R V K E L L - T Y E E P F V V A G D Y N V I P T P E D V --- Y D ...
... P N G N P L G -- T D K F P Y K L A W M D R L R R H A A L R L - A E E Q P F L L L G D Y N V I P E P K D A --- R N ...
... P N G N P L G -- T E K F P Y K L R W M D R L I A H A R T R L - A E E T P F L L L G D Y N V I P E P K D A --- R N ...
... P N G N P L G -- T E K F P Y K L G W M D R L I A H A K R R L - D D E I P Y L L L G D Y N V I P D P M D A --- K N ...
... P N G N P V S A D S V K F P Y K L G W M E R L E A W A Q E R L - E L E E P L I L A G D Y N V I P M P V D C --- H D ...
... P N G N P V D -- T E K F P Y K L R W M E R L Q A F A E D R L - A L E E P L V L A G D Y N V I P E P V D C --- H N ...

completely conserved

PNGENRK-NIEKFTEKKKNFFYRLKKFIEKNF-SPNDVFVLLM**GDINVAINDKD**I**GISENN**...
POGENRN-NKIKFNKKKNFYKNLIK**YVOKKT**-FENKNTTT**MGDMNT**SPEDODIGIDPYH...
PHGESFY-KTDKFEEKKFFYQKLYLF**Y**...
PNGESKK-NLKKFEI**KKIFYKSLFLF**L**NKFRY**-EKNM**HILLIM**GD**MNIS**PED**LDV**G**LSLTS**...
PQGKSID--HPDYQAKHRFFDRLLNLFQREF-SPHTPLLWVG**DMNVAPTDIDV**---TS...
PNGNPAP--GPKFDYKLRWF**ERLRLRAQELI**-ATGAPVVIAG**DYNVMPTEILDV**---YK...
PNGNPAP--GPKFDYKLRWF**DRLITHAQGLL**-AAGKPVLLTGDFN**VMPTEILDV**---YK...
PNGNPVE--TPKYPYKLRWM**DRLLIRYAEDRL**-ALEEPLVLAG**DYNVLPTPDDV**---AN...
PNGNPPQ--TEKYPYKLKWMDRLLAYSKERL-KSEEPFVLAG**DFNVIPTPEDV**---YN...
PNGNPPN--TEKYPYKLKWMSRLRDYARERL-KTEEPLILAG**DFNVIPAAADV**---SN...
PNGNPVG--SEKYPYKLSWMARLRDYAQ**QQL**R-A**EPLVAGDFNVIPQAEDV**---HN...
PNGNPVP--GPKYDYKLA**WMERLRARAI**ELL-KSEAPFVMAG**DYNIIIPQPMDA**---AK...
PNGNPAP--GPKYDYKLA**WMARMHARVES**LL-PL**EPLVFCGDYNVIPQAEDA**---AK...
PNGNPAP--GPKYDYKLA**WMERLEARARE**LL-AEEMPALMA**GDFNVIPQAEDA**---AR...
PNGNPAP--GPKYDYKLA**WMERLRA**REALL-KAEEPALMA**GDFNVIPQAEDA**---AK...
PNGNPAP--GPKFDYKLA**WMQRLEARAK**ALL-AD**EMPFLMAGDYNIIIPQAEDA**---AK...
PNGNPVD--TEKFPYKLEWMDR**LIARAK**ELL-L**LEEPFVMMGDYNIIIPHEDDV**---HD...
PNGNPID--SDKFPYKLSW**MERLRSRV**KELL-TYEEPFVVA**GDFNVIPTPEDV**---YD...
PNGNPLG--TDKFPYKLA**WMDRLRRHA**ALRL-AEEQPFL**LLGDYNVIPEPKDA**---RN...
PNGNPLG--TEKFPYKLRWMDR**LIAHART**RL-AEETPFL**LLGDYNVIPEPKDA**---RN...
PNGNPLG--TEKFPYKL**GWMDR**LIAHAKRRL-DDEIPY**LLLGDYNVIPDPMDA**---KN...
PNGNPVSADSVKFPYKLGW**MERLEAWA**QERL-E**LEEPLILA**GD**DYNVIPMPVDC**---HD...
PNGNPVD--TEKFPYKLRW**MERLQAF**AEDRL-A**LEEPLVLAGDYNVIPEPVDC**---HN...

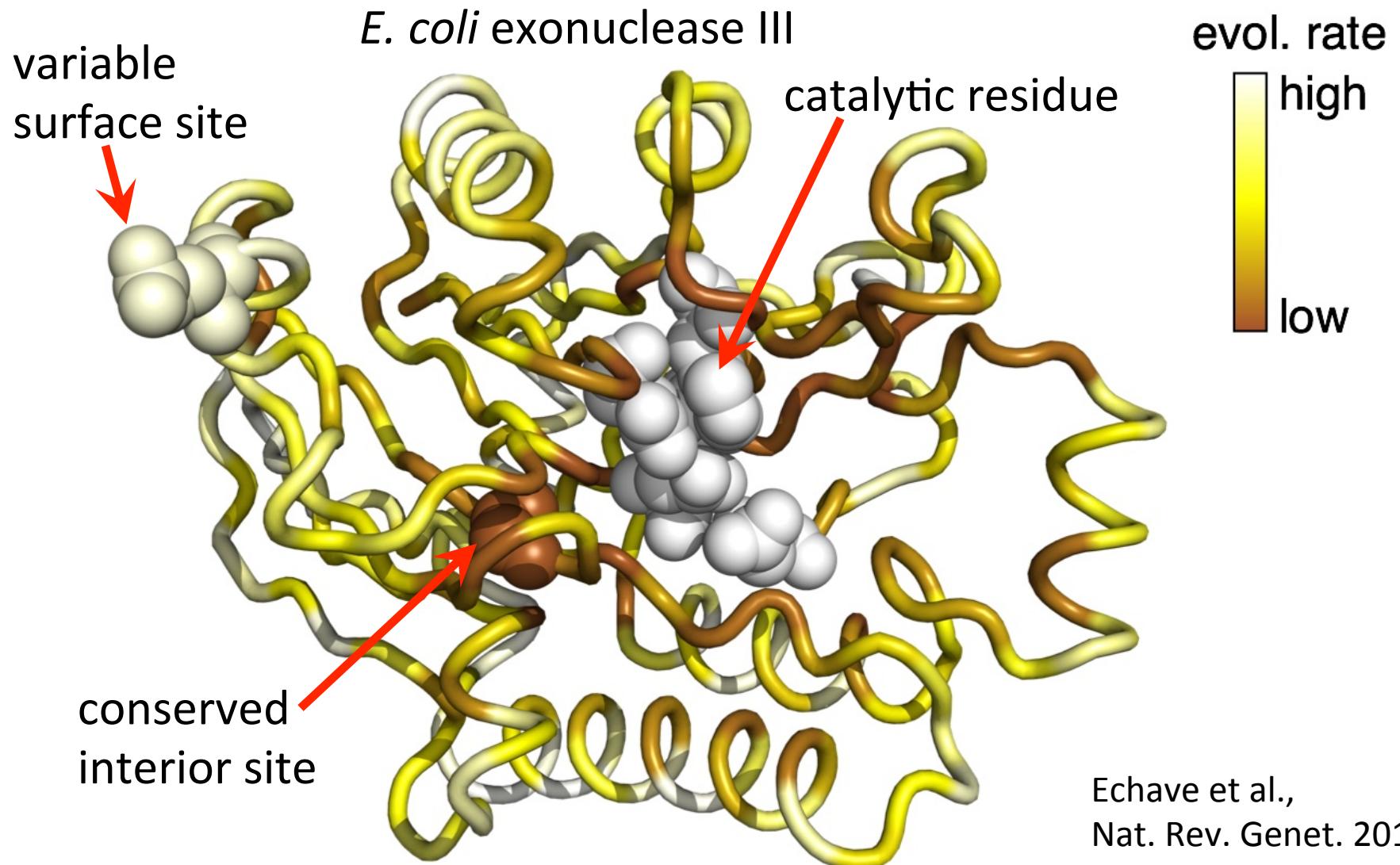
mostly conserved

... P N G E N R K - N I E K F T E K K N F F Y R L K K F I E K N F - S P N D F V I L M G D I N V A I N D K D I G I S E N N ...
... P Q G E N R N - N K I K F N K K K N F Y K N L I K Y V O K K T - F E N K N T T I M G D M N T S P E D Q D I G I D P Y H ...
... P H G E S F Y - K T D K F E E K K F F Y Q K L Y L F - E K N M H I L I M G D M N I S P E D L D V G I L S V N S ...
... P N G E S K K - N L K K F E I K K I F Y K S L F L F L N K F Y - E K N M H I L I M G D M N I S P E D L D V G I L S L T S ...
... P Q G K S I D -- H P D Y Q A K H R F F D R L L N L F Q R E F - S P H T P L L W V G D M N V A P T D I D V - - - - T S ...
... P N G N P A P -- G P K F D Y K L R W F E R L R L R A Q E L I - A T G A P V V I A G D Y N V M P T E L D V - - - - Y K ...
... P N G N P A P -- G P K F D Y K L R W F D R L I T H A Q G L L - A A G K P V L L T G D F N V M P T E L D V - - - - Y K ...
... P N G N P V E -- T P K Y P Y K L R W M D R L I R Y A E D R L - A L E E P L V L A G D Y N V L P T P D D V - - - - A N ...
... P N G N P P Q -- T E K Y P Y K L K W M D R L L A Y S K E R L - K S E E P F V L A G D F N V I P T P E D V - - - - Y N ...
... P N G N P P N -- T E K Y P Y K L K W M S R L R D Y A R E R L - K T E E P L I L A G D F N V I P A A A D V - - - - S N ...
... P N G N P V G -- S E K Y P Y K L S W M A R L R D Y A Q Q R L - K T E E P L I L A G D F N V I P Q A E D V - - - - H N ...
... P N G N P V P -- G P K Y D Y K L A W M E R L R A R A I E L L - K S E A P F V M A G D Y N I I P Q P M D A - - - - A K ...
... P N G N P A P -- G P K Y D Y K L A W M A R M H A R V E S L L - P L E E P L V F C G D Y N V I P Q A E D A - - - - A K ...
... P N G N P A P -- G P K Y D Y K L A W M E R L E A R A R E E L L - A E E M P A L M A G D Y N V I P Q A E D A - - - - A R ...
... P N G N P A P -- G P K Y D Y K L A W M E R L R A R A E A L L - K A E E P A L M A G D Y N V I P Q A E D A - - - - A K ...
... P N G N P A P -- G P K F D Y K L A W M Q R L E A R A K A L L - A D E M P F L M A G D Y N I I P Q A E D A - - - - A K ...
... P N G N P V D -- T E K F S Y K L E W M D R L I A R A K E L L - L L E E P F V M M G D Y N I I P H E D D V - - - - H D ...
... P N G N P I D -- S D K F P Y K L S W M E R L R S R V K E L L - T Y E E P F V V A G D Y N V I P T P E D V - - - - Y D ...
... P N G N P L G -- T D K F P Y K L A W M D R L R R H A A L R L - A E E Q P F L L L G D Y N V I P E P K D A - - - - R N ...
... P N G N P L G -- T E K F P Y K L R W M D R L I A H A R T R L - A E E T P F L L L G D Y N V I P E P K D A - - - - R N ...
... P N G N P L G -- T E K F P Y K L G W M D R L I A H A K R R L - D D E I P Y L L L G D Y N V I P D P M D A - - - - K N ...
... P N G N P V S A D S V K F P Y K L G W M E R L E A W A Q E R L - E L E E P L I L A G D Y N V I P M P V D C - - - - H D ...
... P N G N P V D -- T E K F P Y K L R W M E R L Q A F A E D R L - A L E E P L V L A G D Y N V I P E P V D C - - - - H N ...

quite variable

... PNGENRK - NIEKFTEKKNFFYRLKKFIEKNF - SPNDVFVLLMGDINVAINDKD**I**GISENN ...
... PQGENRN - NKIKFNKKKNFYKNL**I**IKYVOKKT - FENKNTT**I**M**I**DMNISPEDQDIGIDPYH ...
... PHGESFY - KTDKFEEKKFFYQKL**I**YLF - E**I**N**I**DMNISPTDLDIGLSVNS ...
... PNGESKK - NLKKFEI**I**KKIFYKSLFLFL**I**Y - E**I**N**I**DMNISPEDLDVG**I**LSLTS ...
... PQGKSID - HPDYQAKHRFFDR**I**LLNL**I**FQREF - SPHTPLLWVG**I**DMNVAPTDIDV - - TS ...
... PNGNPAP - GPKFDYKLRWF**I**ERL**I**RLRAQELI - ATGAPVVIAGDYNVM**I**PTELDV - - YK ...
... PNGNPAP - GPKFDYKLRWF**I**DR**I**LTHAQGLL - AAGKPVL**I**LTGDFNVM**I**PTELDV - - YK ...
... PNGNPVE - TPKYPYKLRWM**I**DRL**I**RYAEDRL - ALEEPLVLAGDYNVLPTPDDV - - AN ...
... PNGNPPQ - TEKYPYKLKWMDR**I**LLAYSKERL - KSEE**I**PFVLAGDFNVIPTPEDV - - YN ...
... PNGNPPN - TEKYPYKLKWMSRL**I**RDYARERL - KTEEPLILAGDFNVIPAAADV - - SN ...
... PNGNPVG - SEKYPYKLSWMARL**I**RDYAQ**I**QL - KTEEPLILAGDFNVIPQAEDV - - HN ...
... PNGNPVP - GPKYDYKLA**I**W**I**MERL**I**RARA**I**ELL - KSEAPFVMAGDYNII**I**PQPM**I**DA - - AK ...
... PNGNPAP - GPKYDYKLA**I**W**I**M**I**ARM**I**HARV**I**ESLL - PLEEPLVFCGDYNVIPQAEDA - - AK ...
... PNGNPAP - GPKYDYKLA**I**W**I**MERL**I**EAR**I**RELL - AEEMPALMAGDYNVIPQAEDA - - AR ...
... PNGNPAP - GPKYDYKLA**I**W**I**MERL**I**RARE**I**ALL - KAEEPALMAGDYNVIPQAEDA - - AK ...
... PNGNPAP - GPKFDYKLA**I**W**I**MQRI**I**EARAK**I**ALL - ADEM**I**MPFLMAGDYNII**I**PQAEDA - - AK ...
... PNGNPVD - TEKFPYKLEWMDR**I**LI**I**ARAK**I**ELL - LLEEPFVMMGDYNII**I**PHEDDV - - HD ...
... PNGNPID - SDKFPYKLSW**I**MERL**I**RSRV**I**KELL - TYEEPFVVAGDYNVIPTPEDV - - YD ...
... PNGNPLG - TD**I**KFPYKLA**I**W**I**MDR**I**RR**I**HA**I**ALRL - AEEQPFL**I**LG**I**DYNVIPEPKDA - - RN ...
... PNGNPLG - TEKFPYKLRWMDR**I**IA**I**HART**I**RL - AEETPF**I**LL**I**GDYNVIPEPKDA - - RN ...
... PNGNPLG - TEKFPYKL**I**GWMDR**I**IA**I**HAK**I**RL - DDEIPY**I**LL**I**GDYNVIPDPMDA - - KN ...
... PNGNPVSAD - SVKFPYKL**I**GW**I**MERL**I**EA**I**WA**I**Q**I**ERL - ELEEPLILAGDYNVIPMPVDC - - HD ...
... PNGNPVD - TEKFPYKLRW**I**MERL**I**Q**I**A**I**FA**I**ED**I**RL - ALEEPLVLAGDYNVIPEPVDC - - HN ...

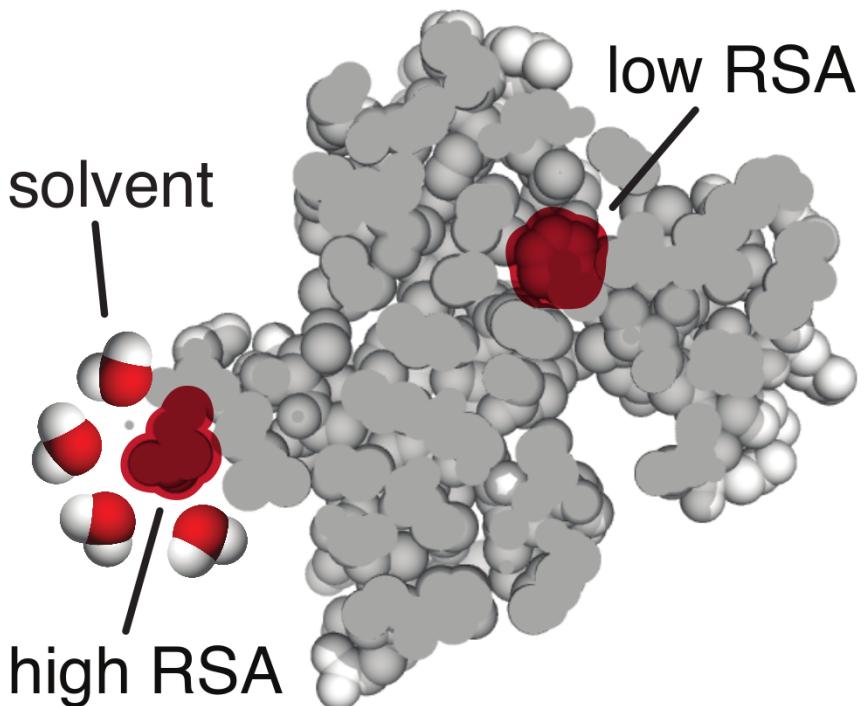
Patterns of sequence variation make more sense in a structural context



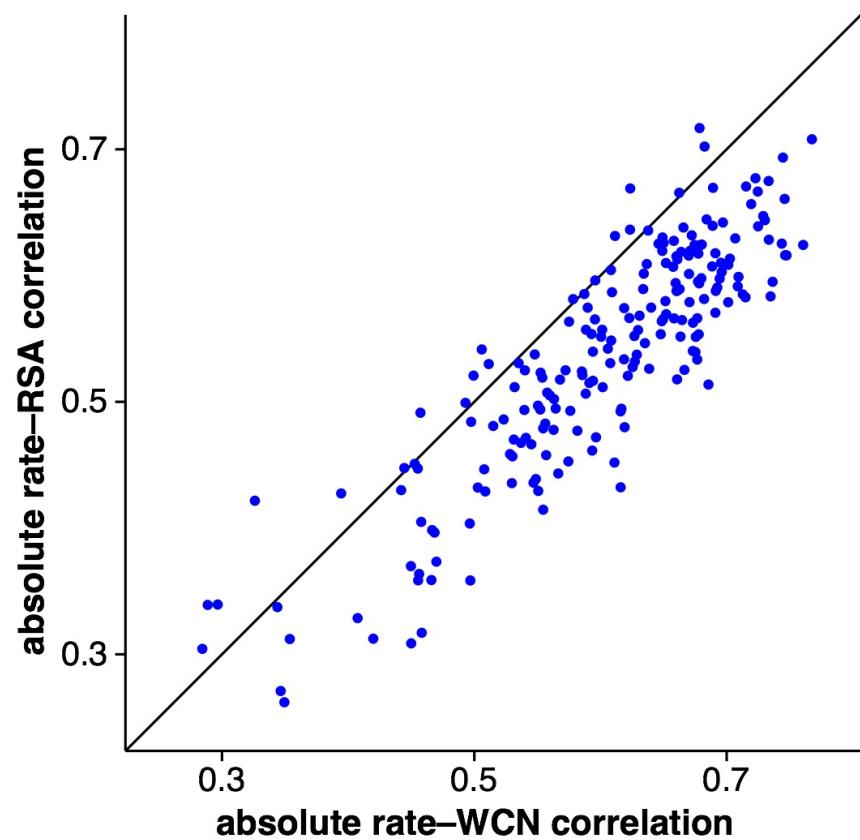
Echave et al.,
Nat. Rev. Genet. 2016

We measure structure with RSA and WCN

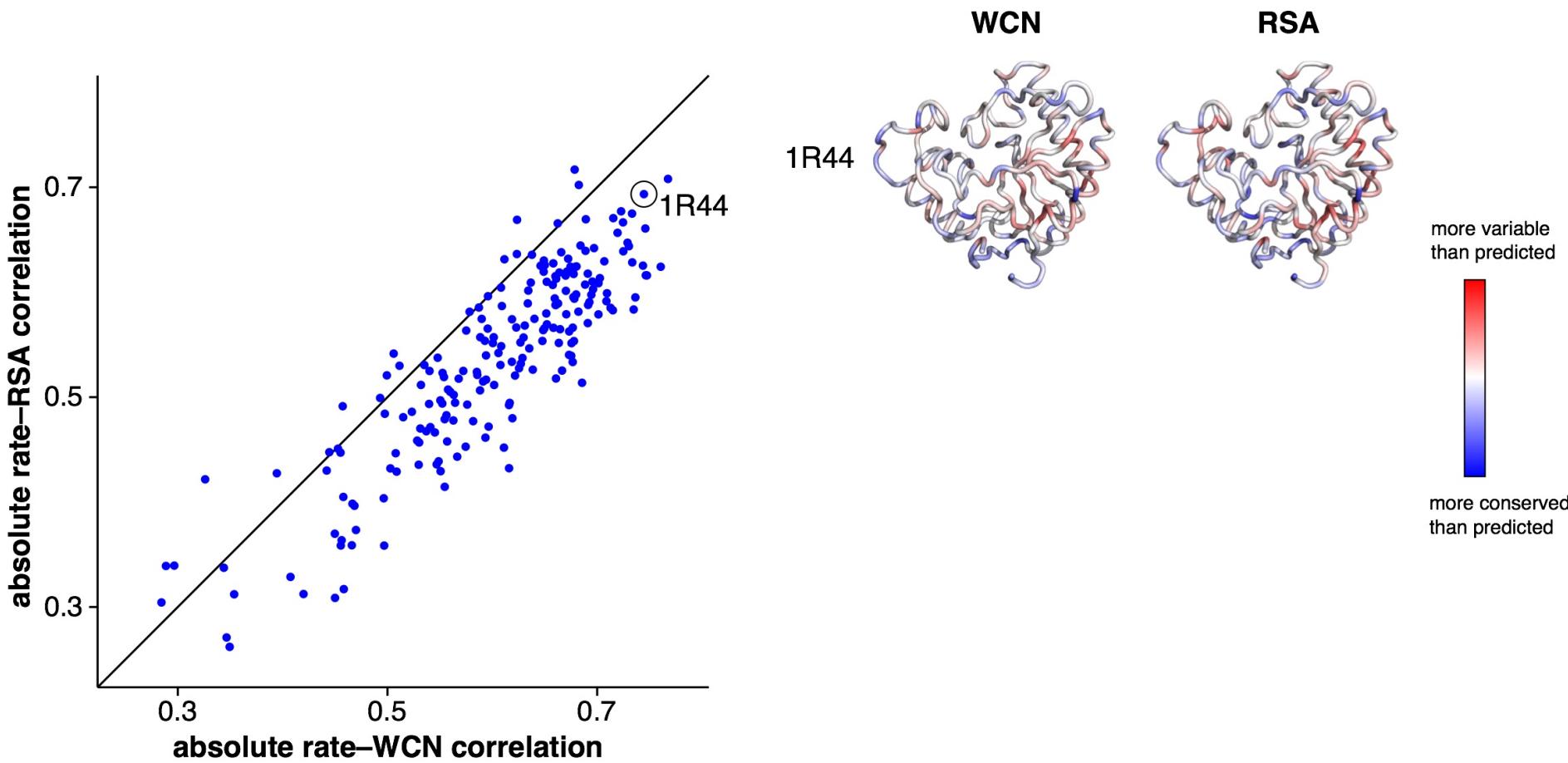
Relative Solvent Accessibility



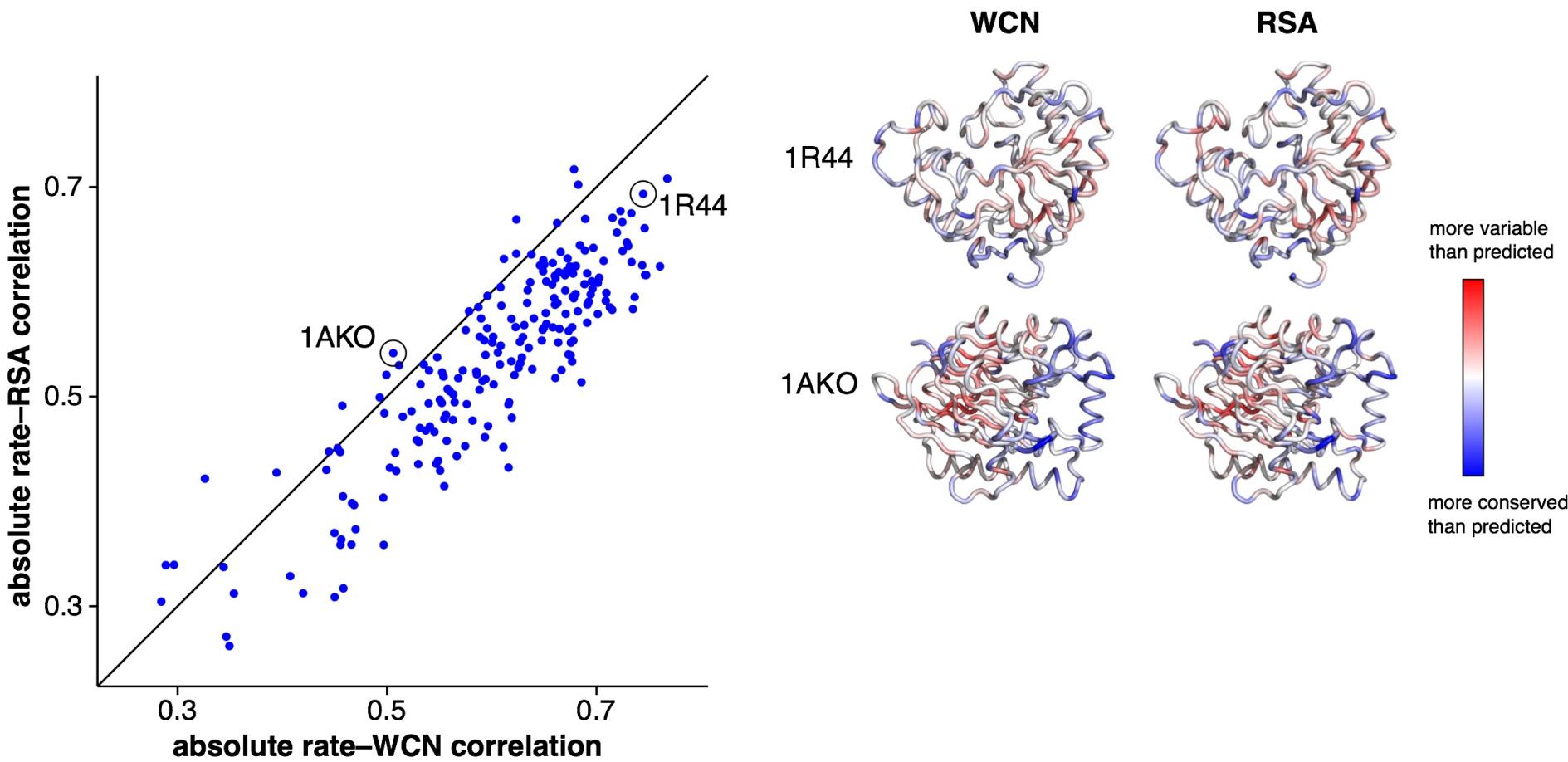
Structure is a good rate predictor for some proteins, but not so much for others



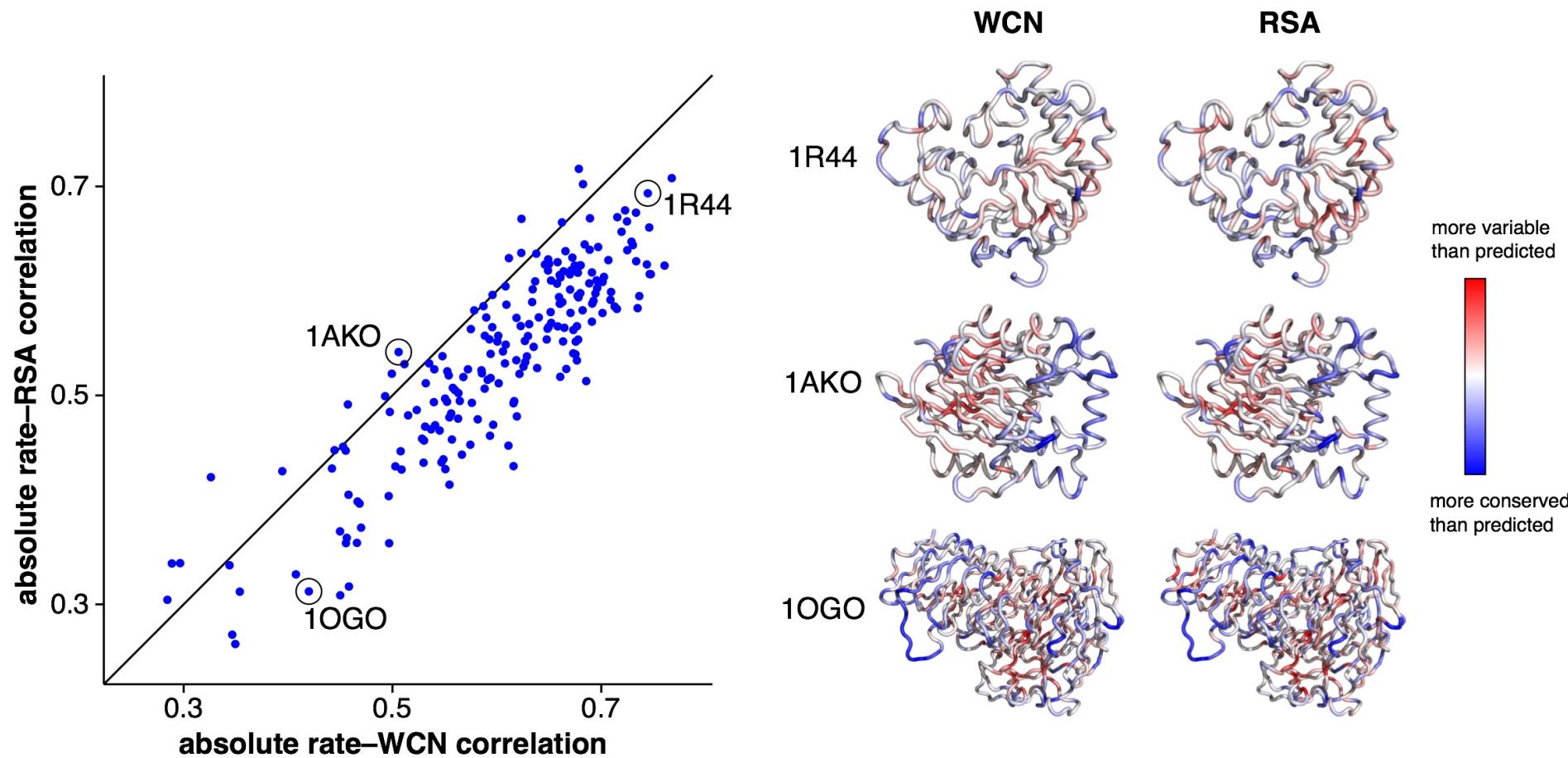
Structure is a good rate predictor for some proteins, but not so much for others



Structure is a good rate predictor for some proteins, but not so much for others



Structure is a good rate predictor for some proteins, but not so much for others

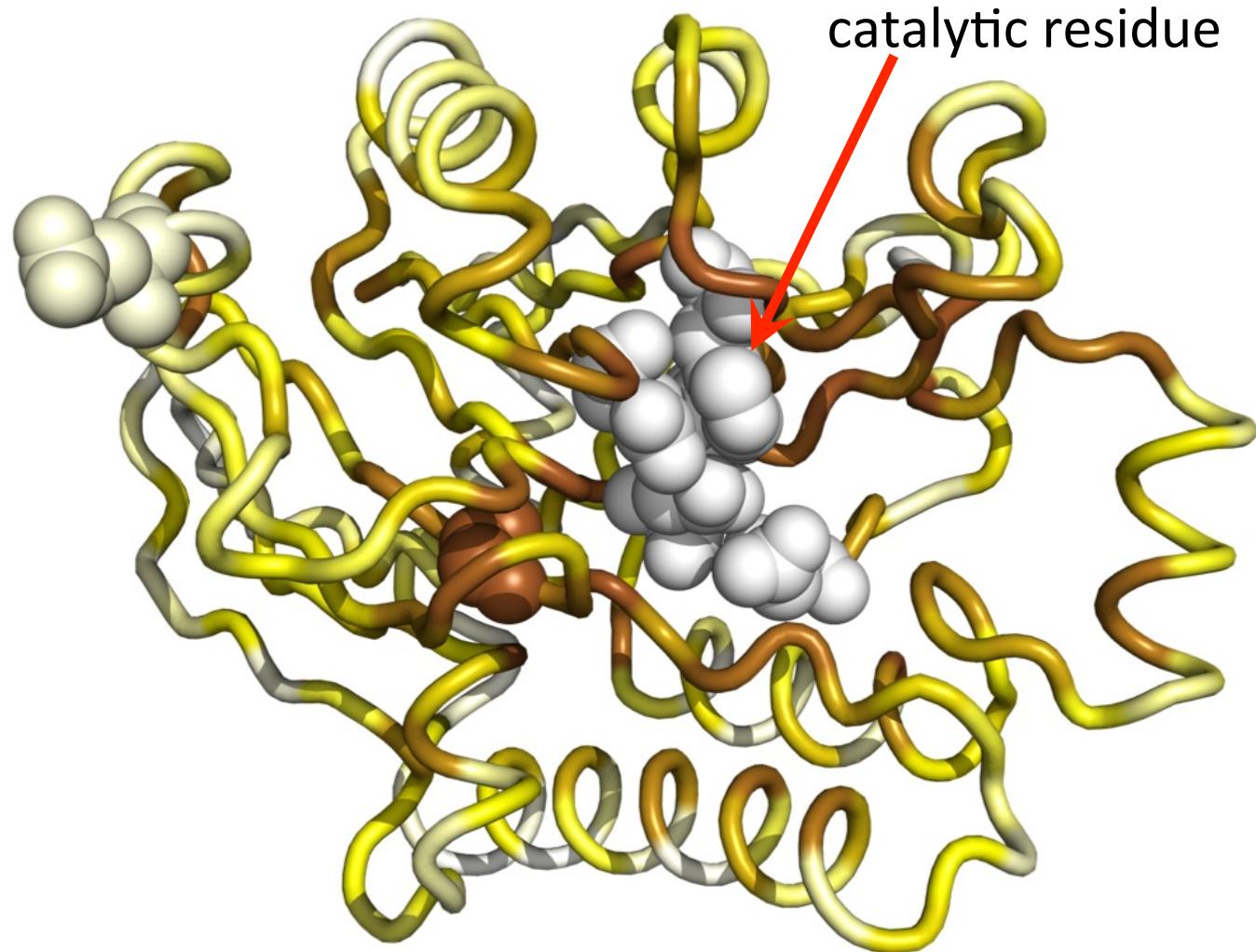


Part I: How does protein function constrain site-specific evolution?

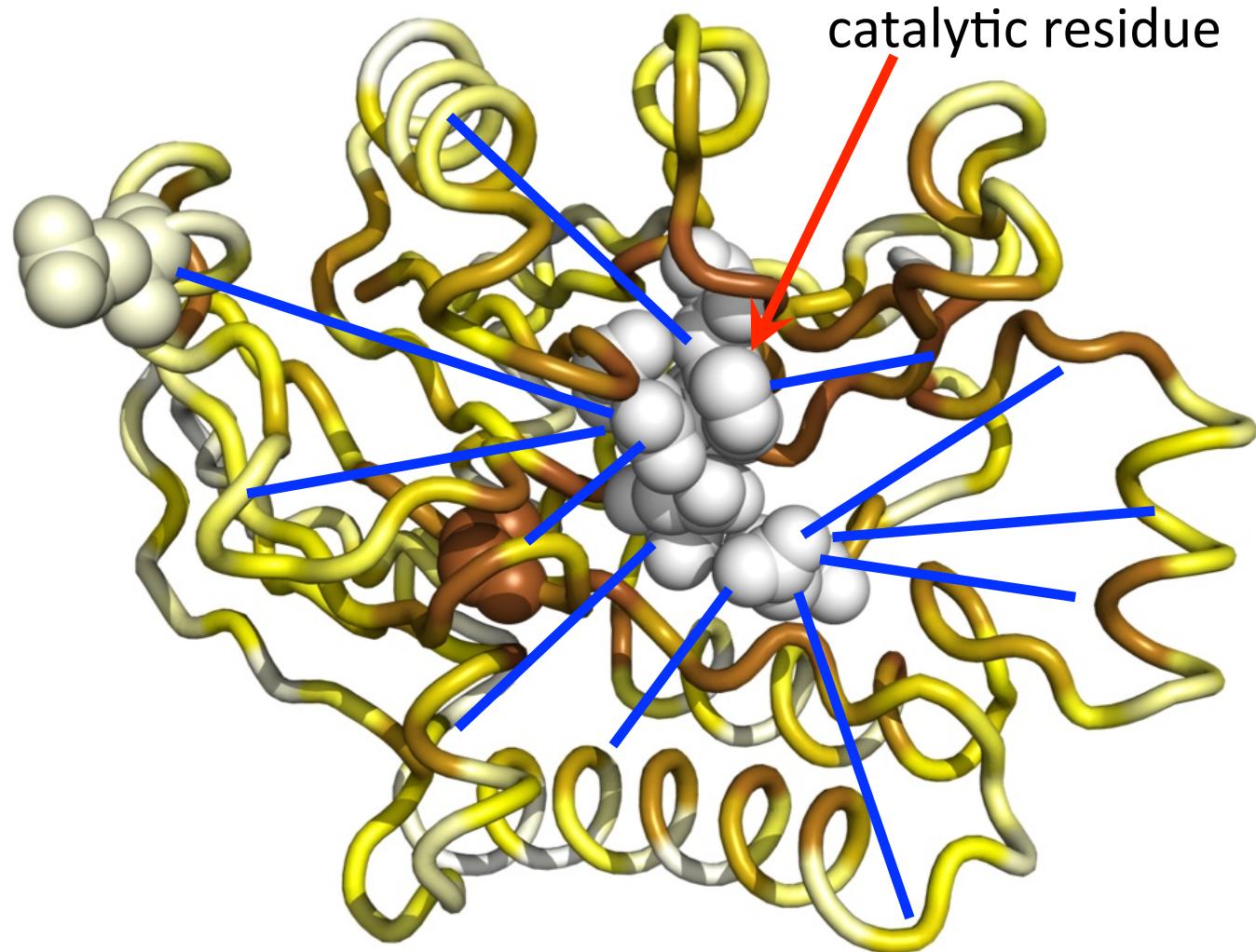
Work by graduate student Ben Jack

In collaboration with Julian Echave

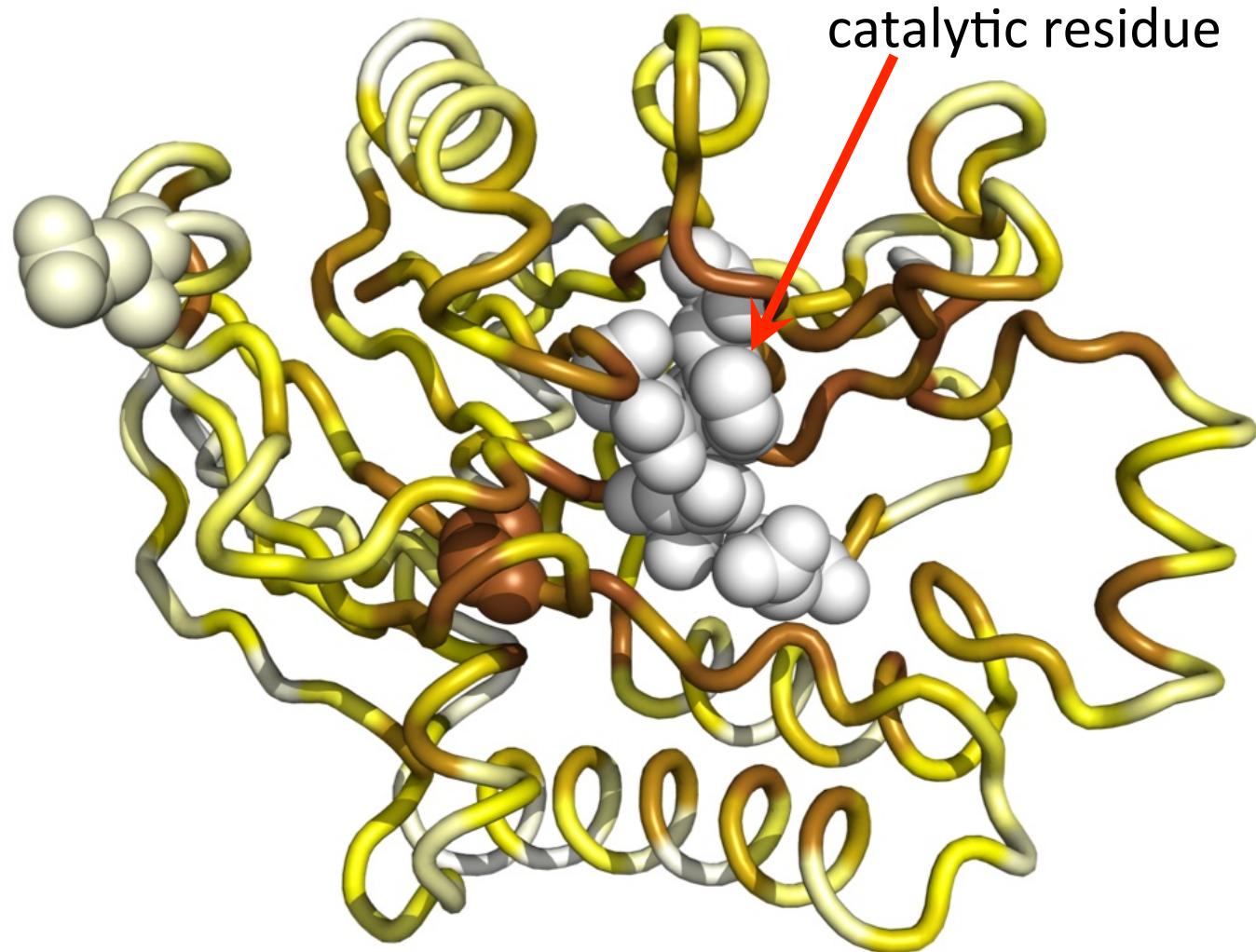
We measure function with distance to the nearest catalytic residue



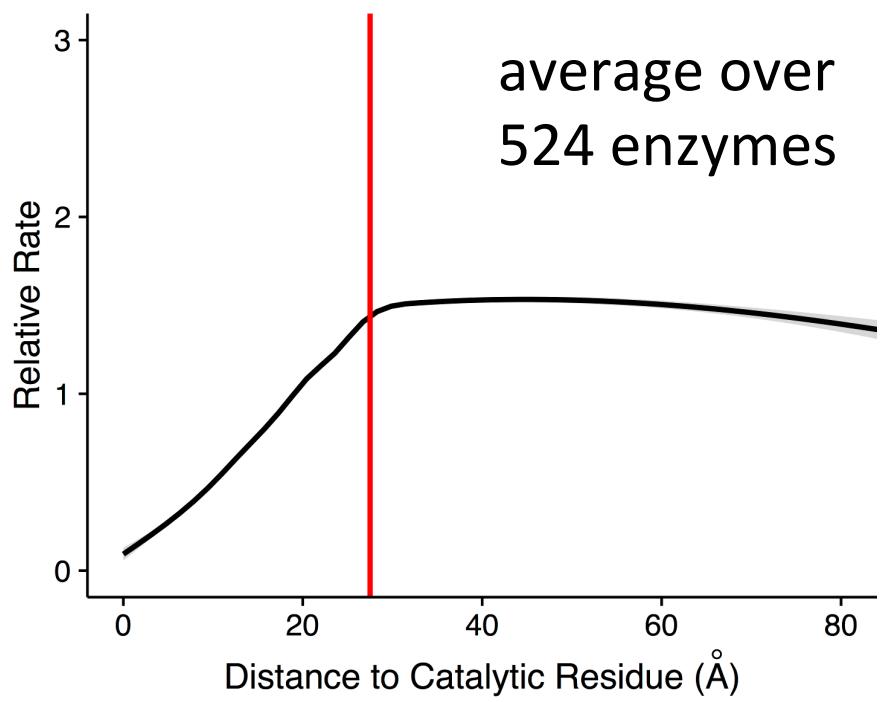
We measure function with distance to the nearest catalytic residue



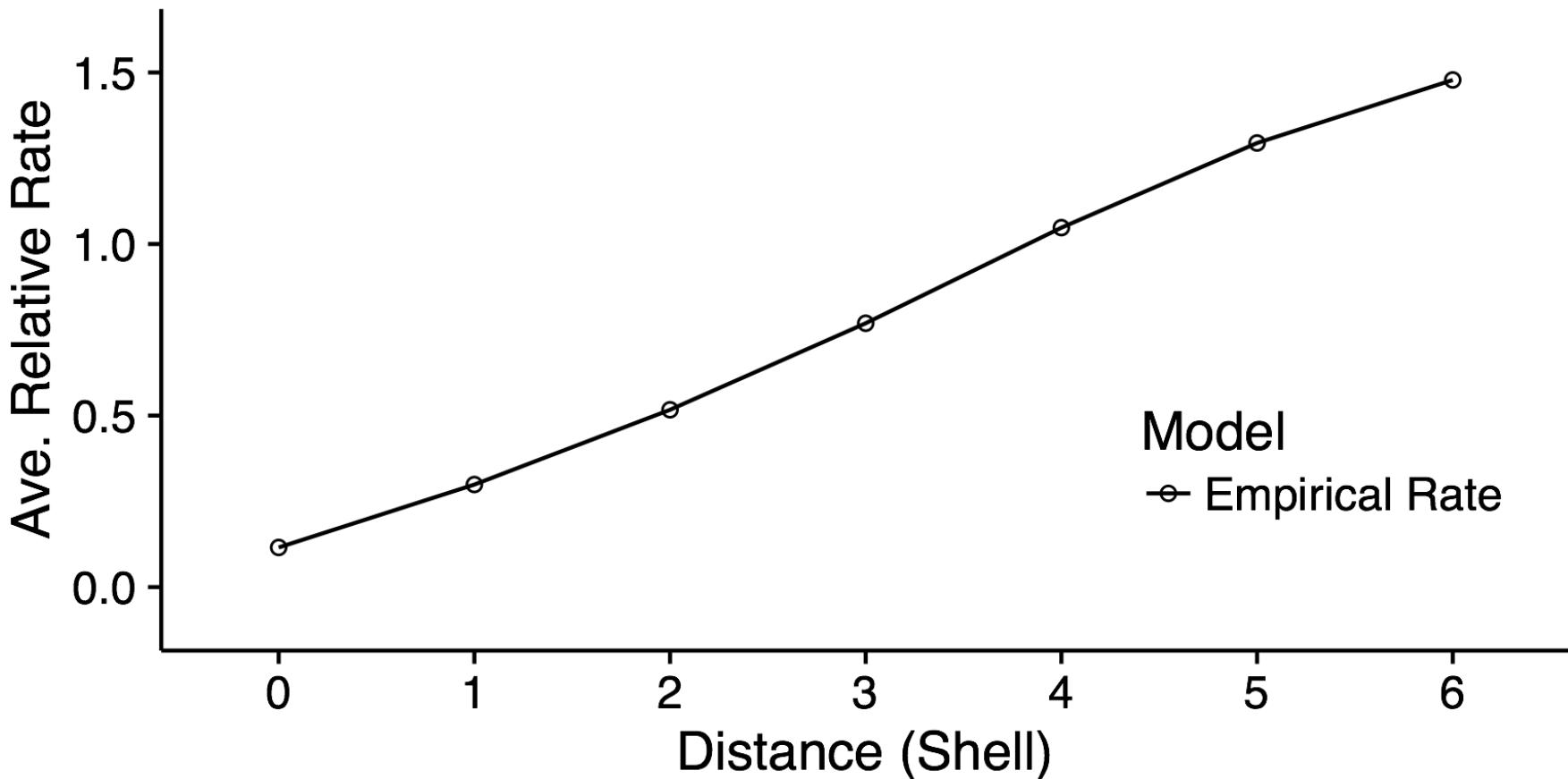
We measure function with distance to the nearest catalytic residue



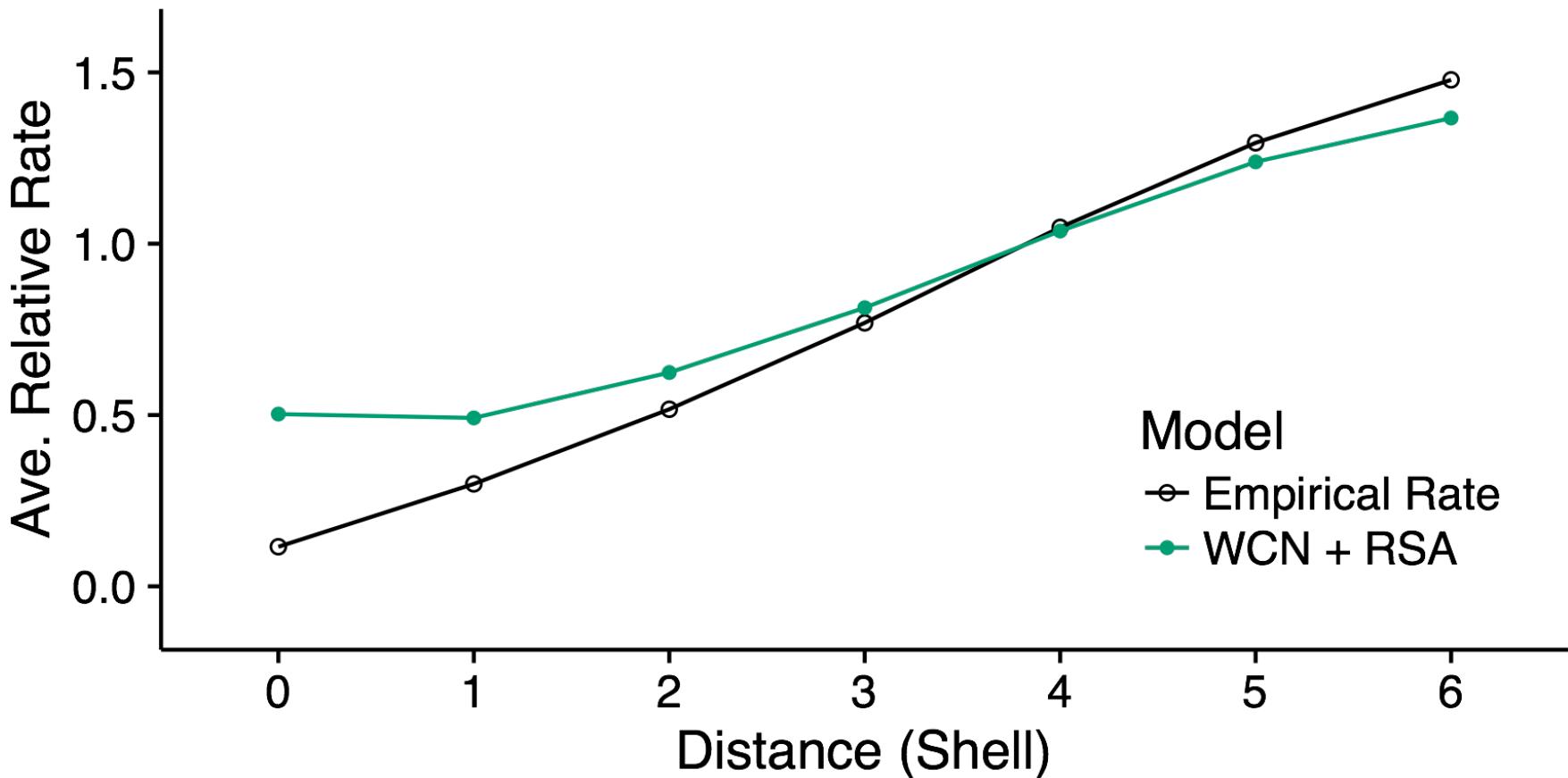
Catalytic residues impose long-range conservation gradients



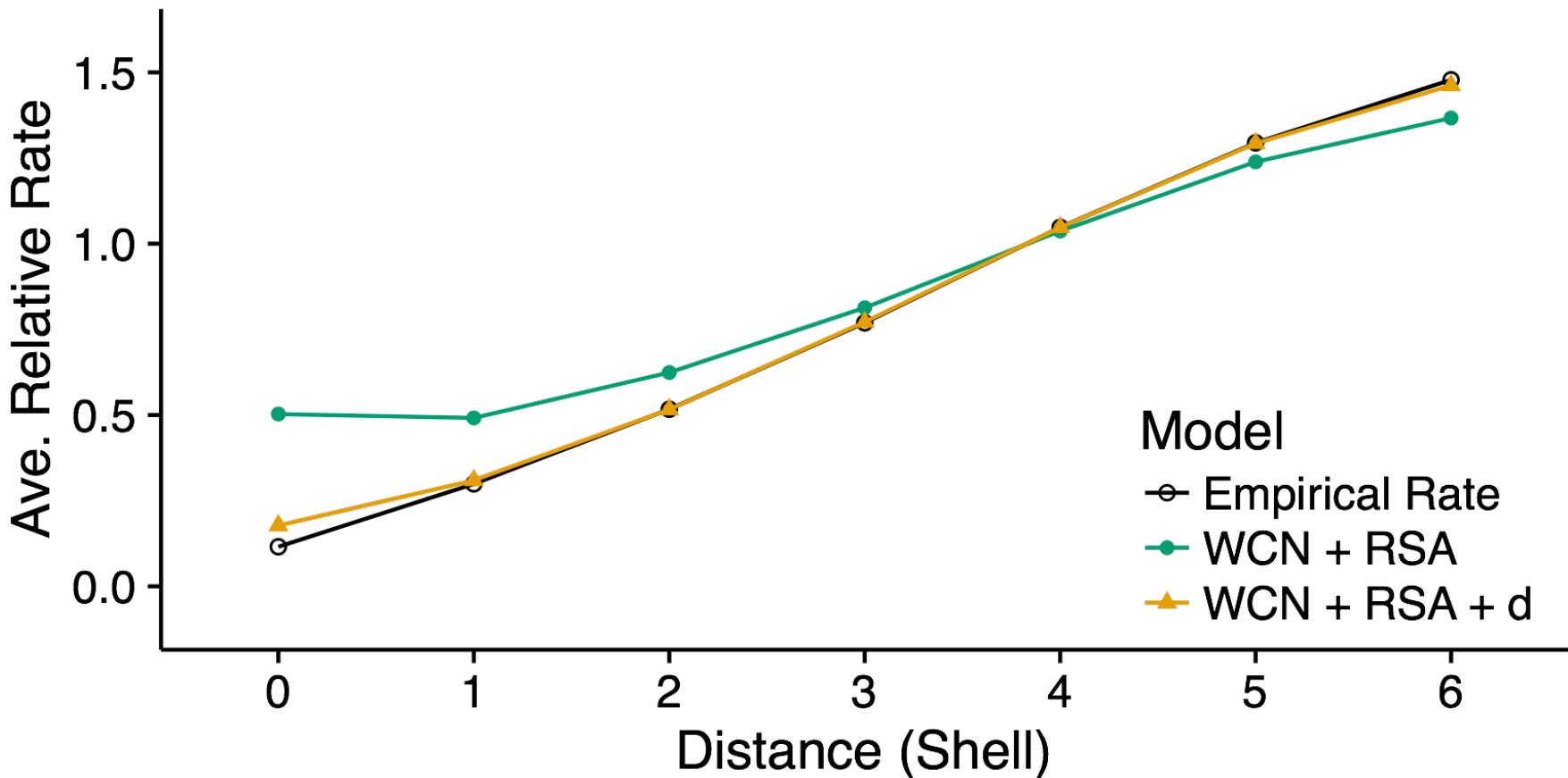
Distance is critical to accurately predict average rate near catalytic residues



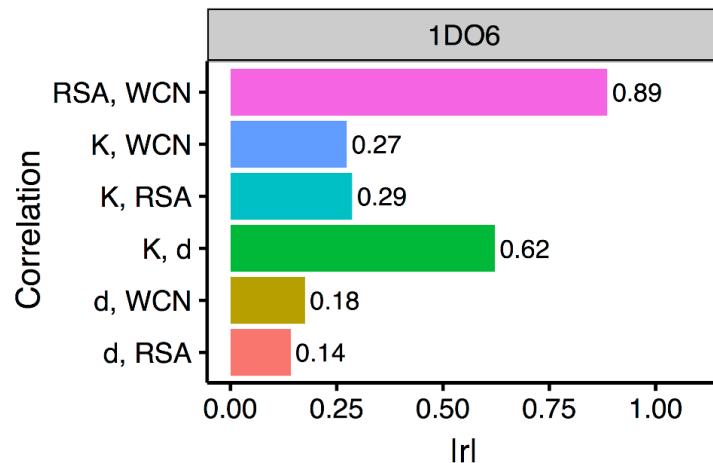
Distance is critical to accurately predict average rate near catalytic residues

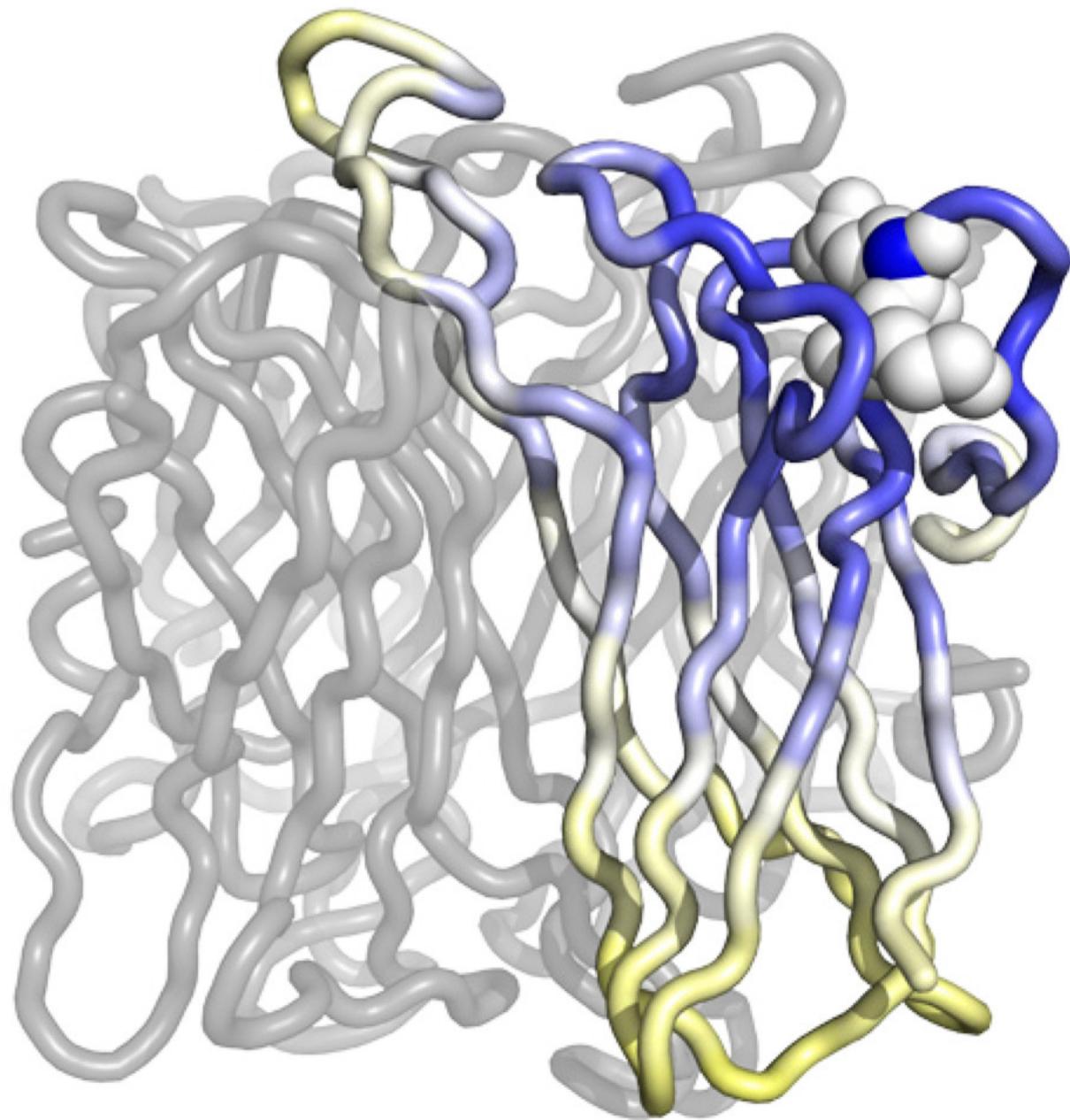


Distance is critical to accurately predict average rate near catalytic residues

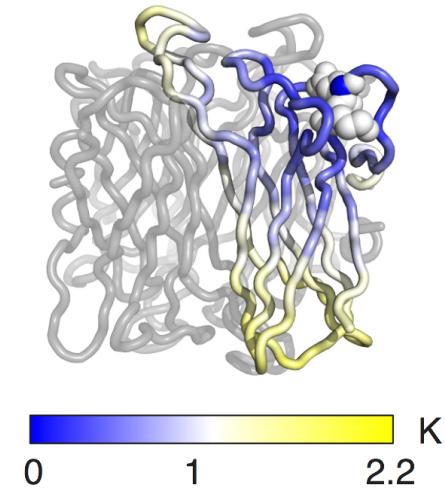
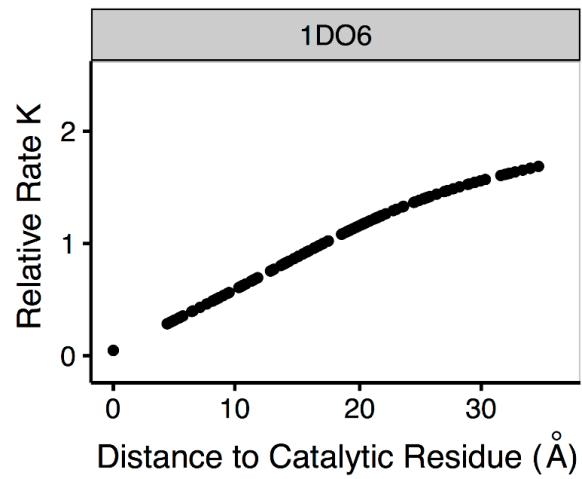
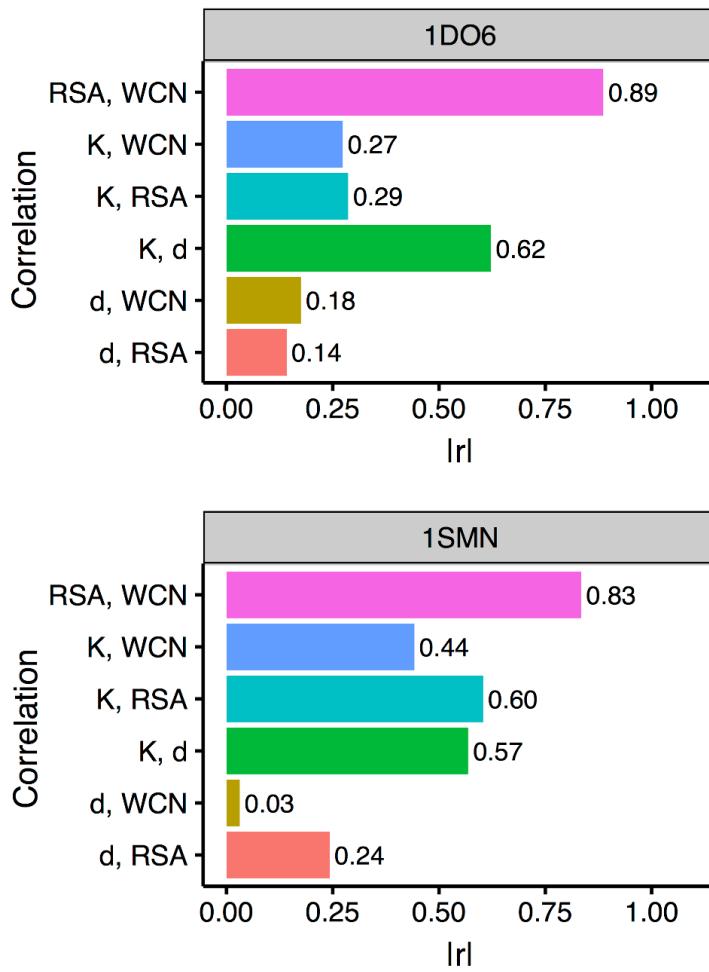


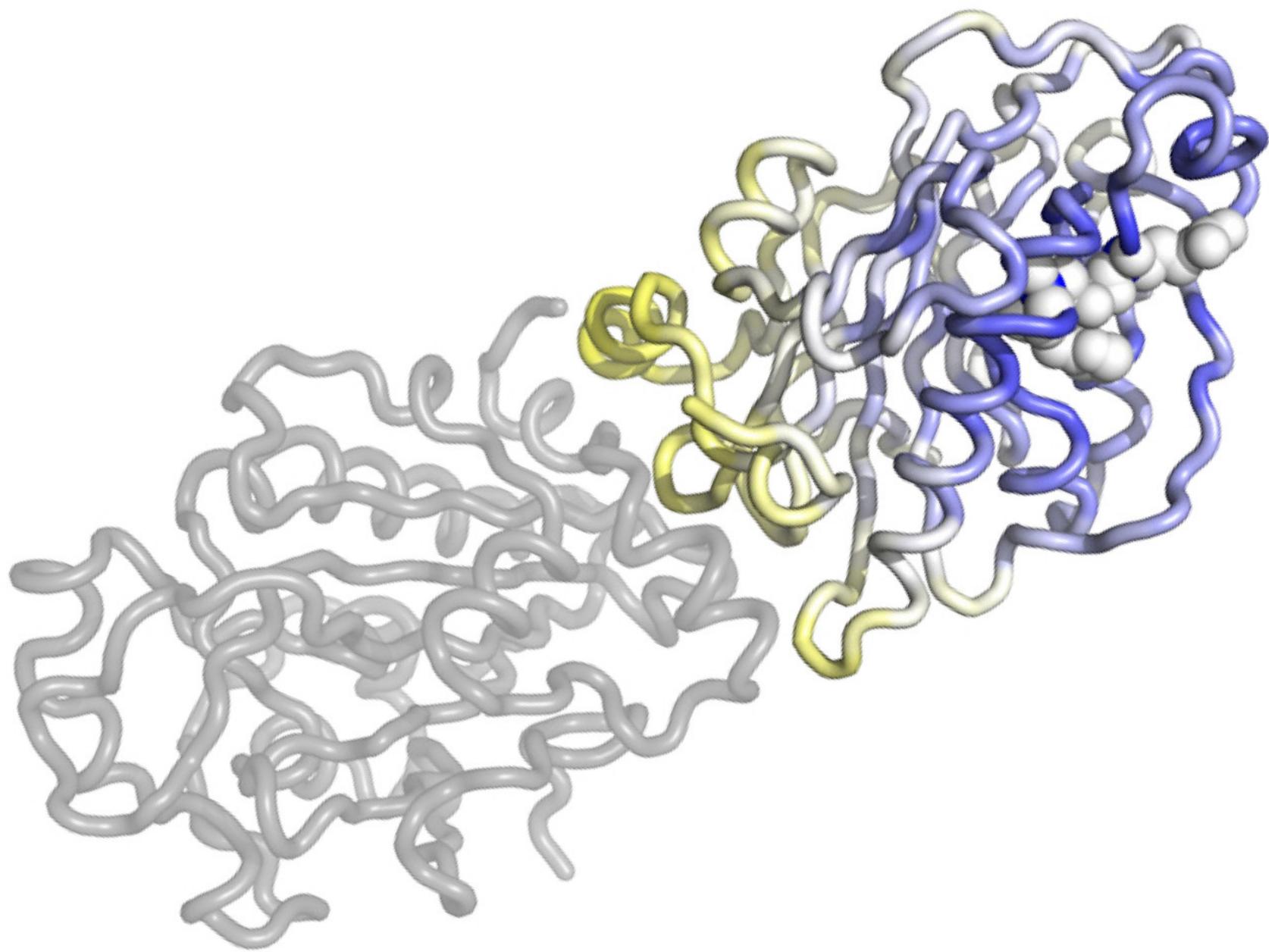
Effect is strongest for active sites on the protein surface



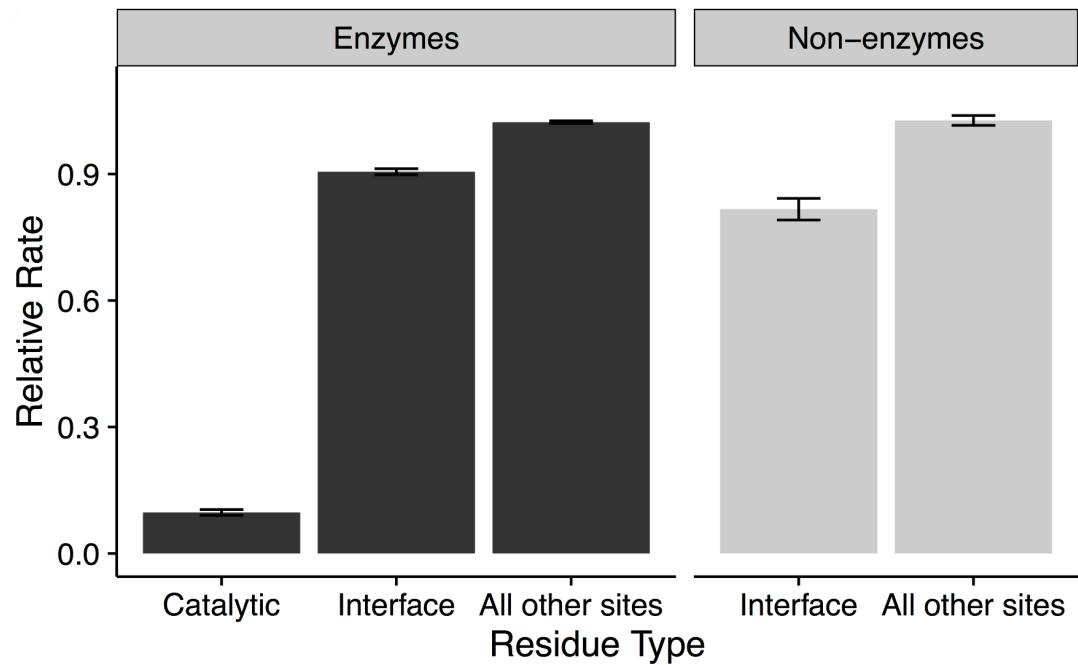


Effect is strongest for active sites on the protein surface





Catalytic residues impose much stronger constraints than interface residues



Take-home message

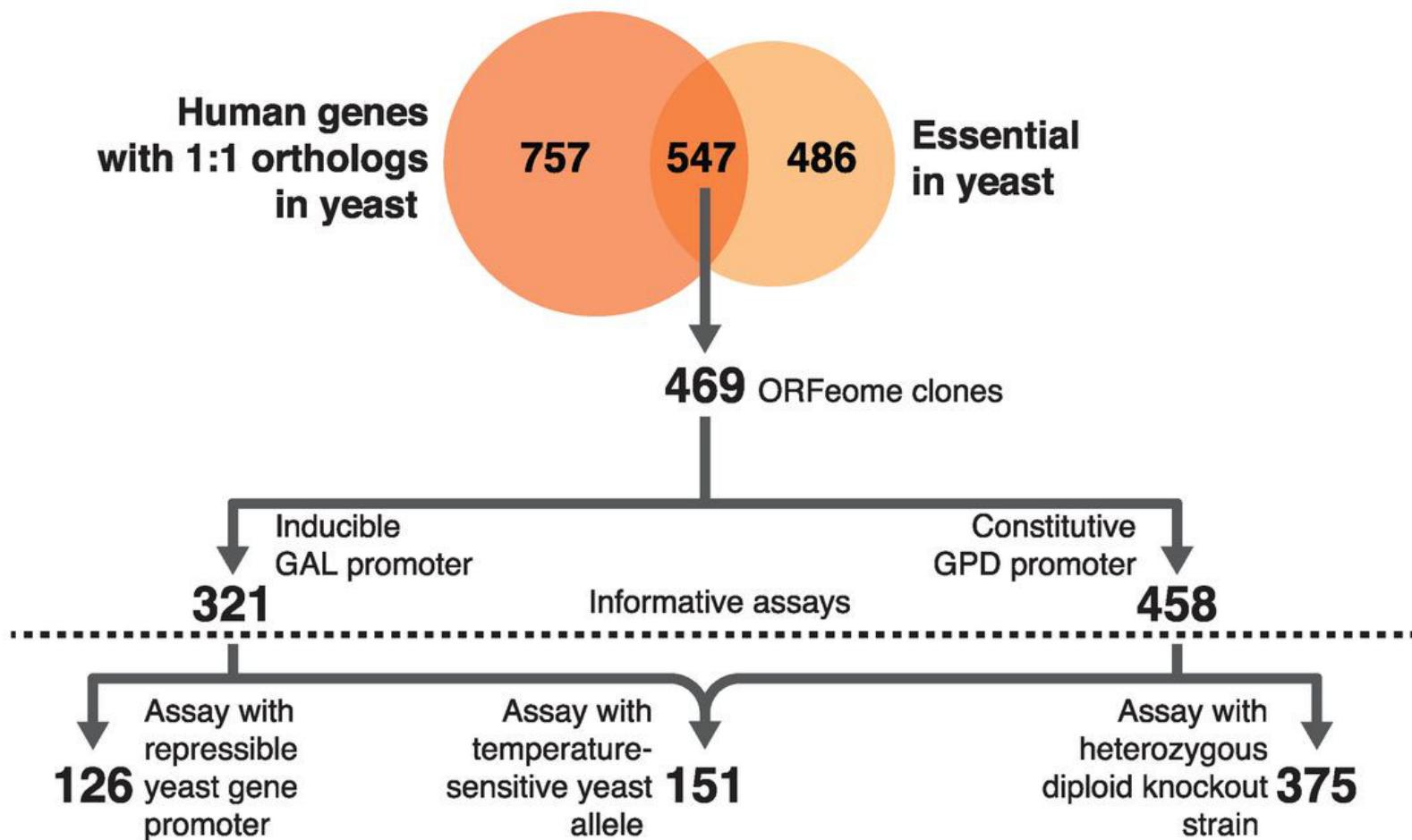
Catalytic residues generate long-range selection gradients covering most of a typical enzyme structure.

Part II: When can proteins interact with ancestral partners?

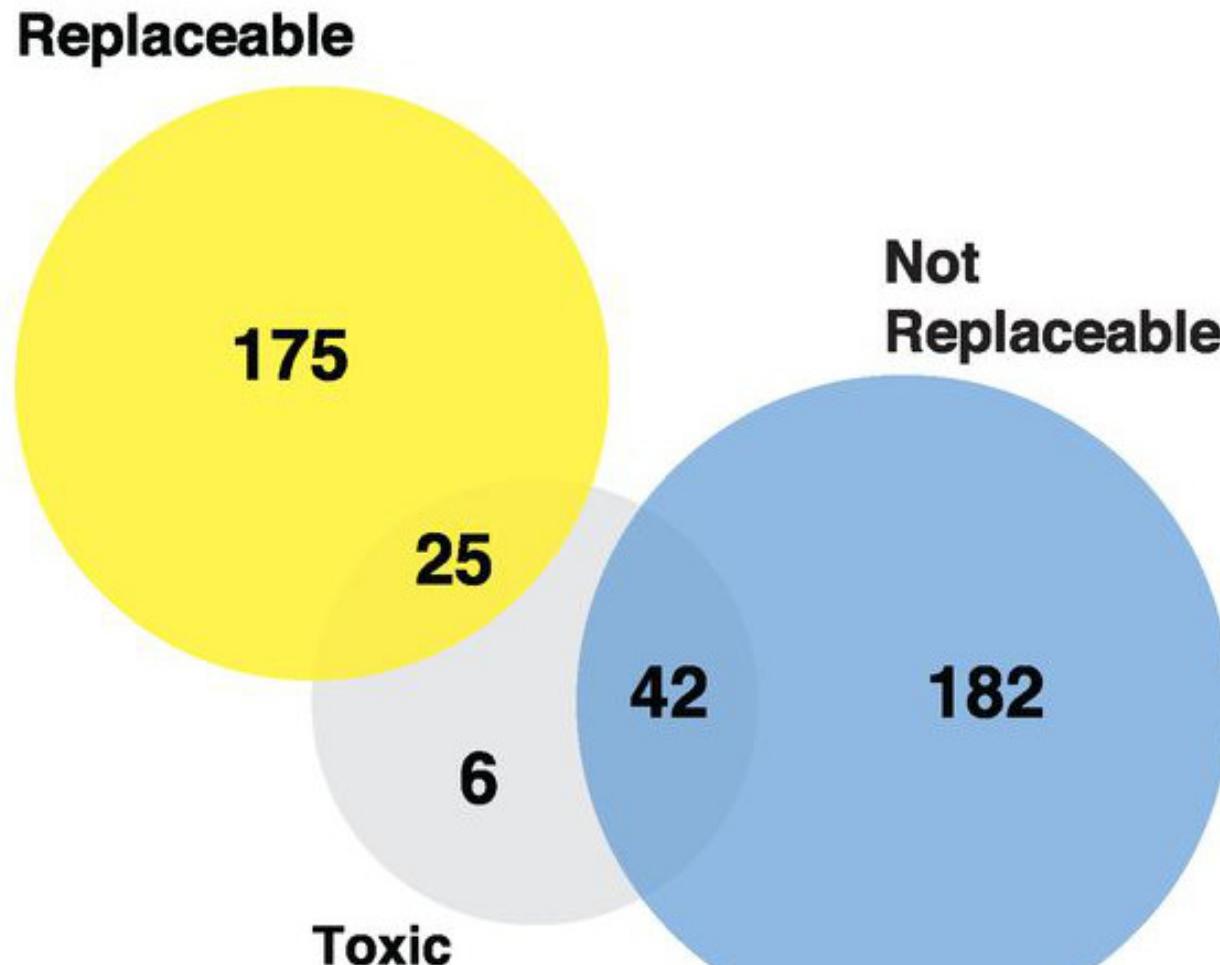
In collaboration with the Marcotte lab

Protein evolution simulations by Austin Meyer

Can human genes function in a distant relative (baker's yeast)?

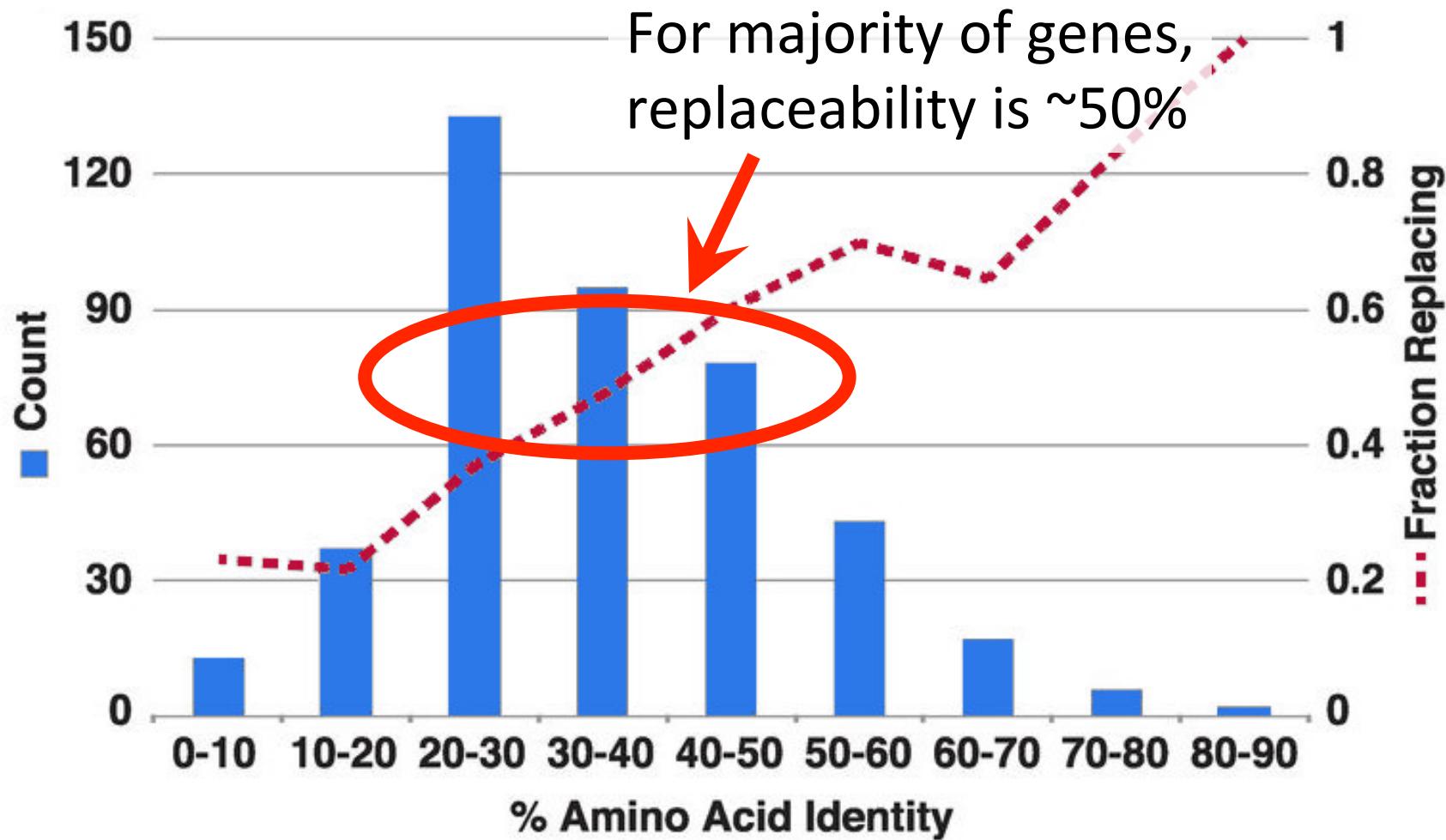


Yes!

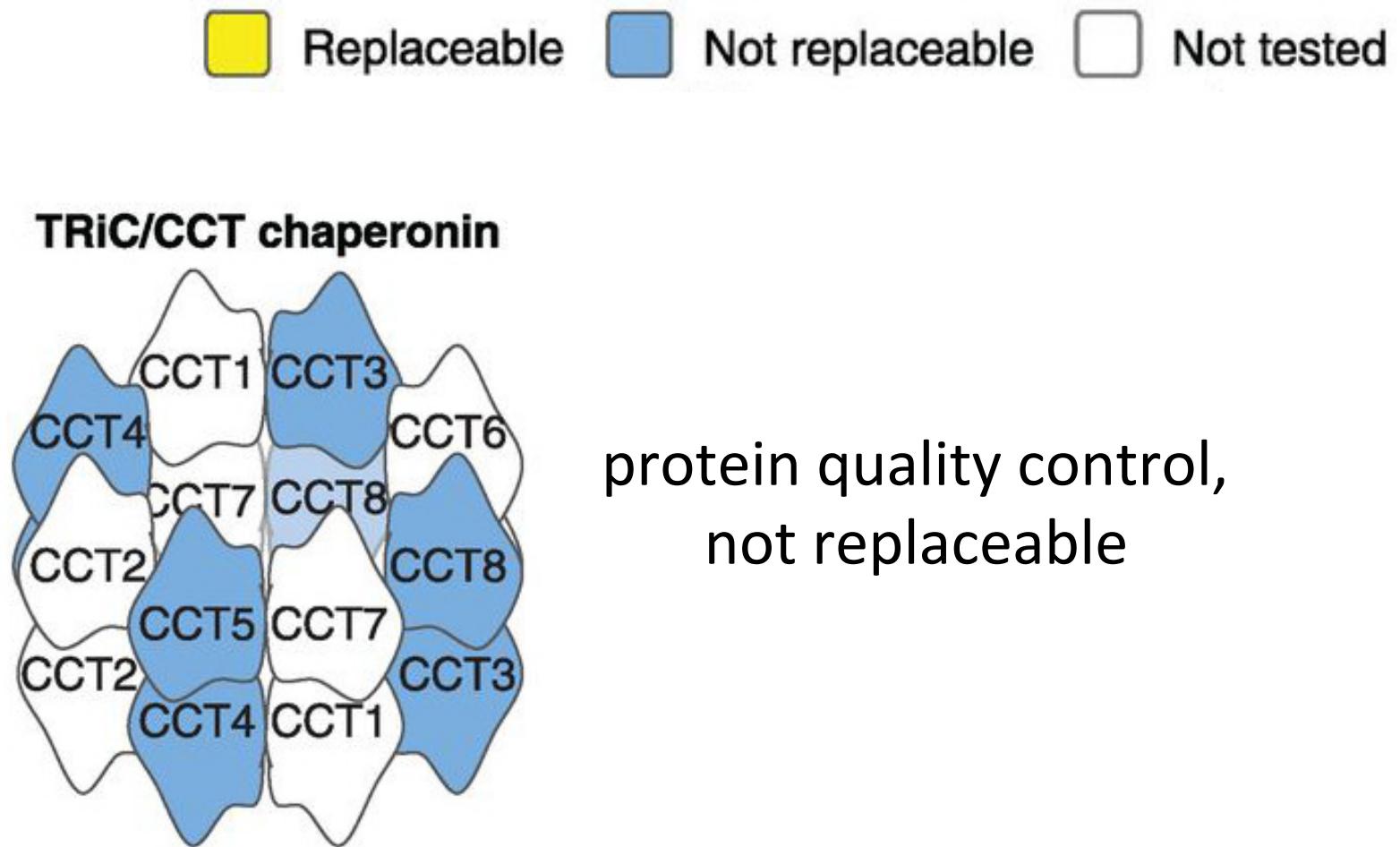


What predicts whether a gene will
be replaceable?

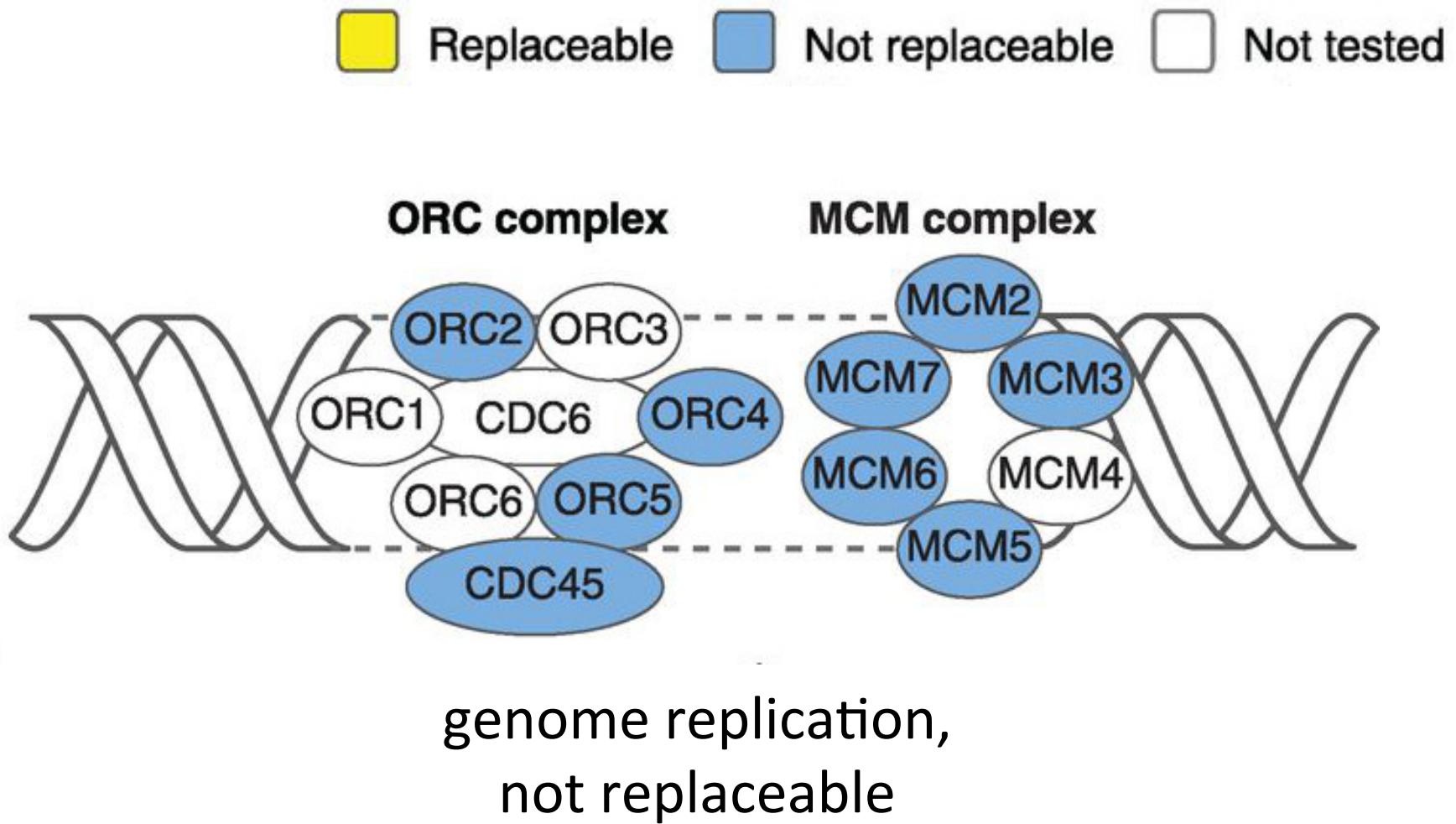
Sequence divergence is not a good predictor of replaceability



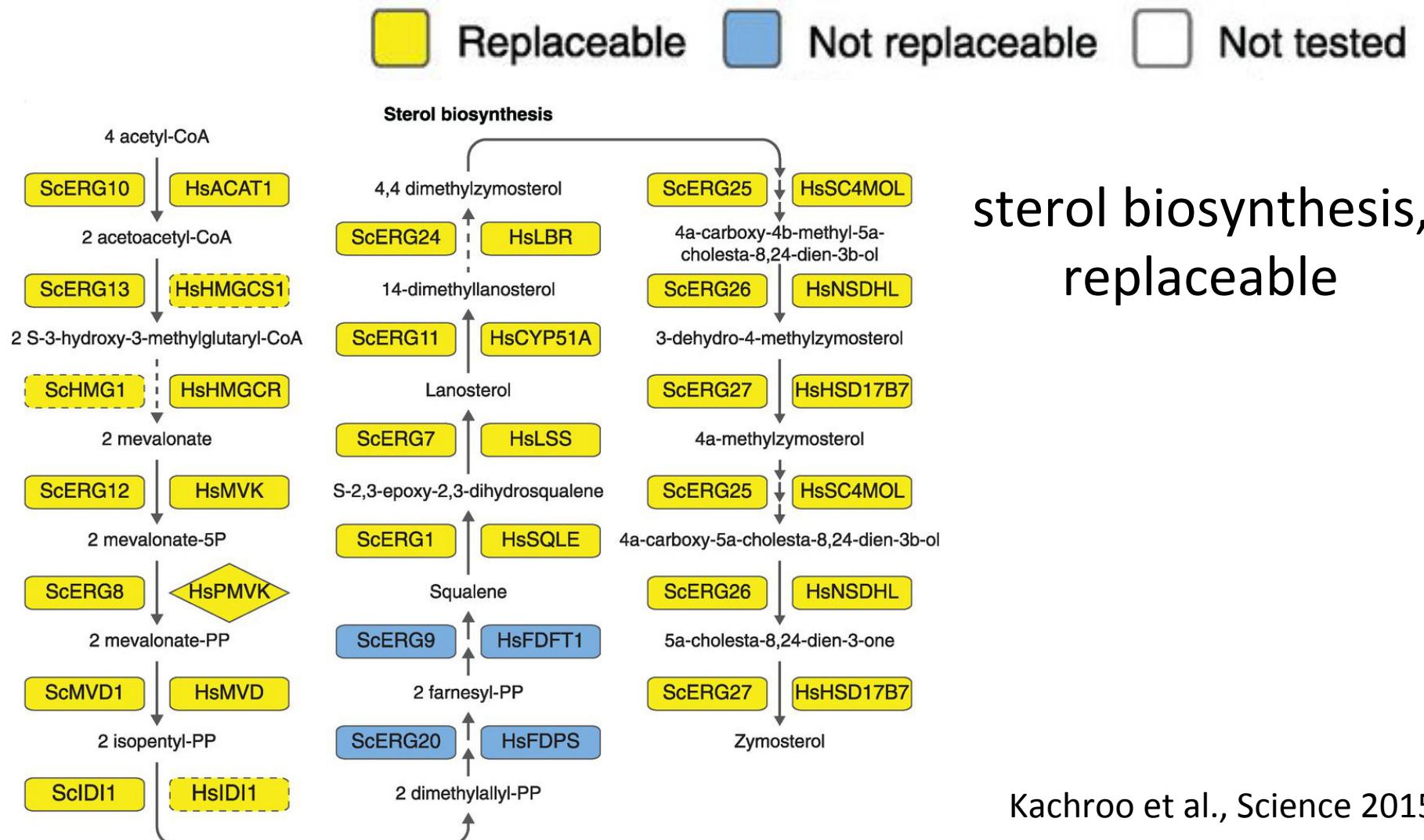
Genes tend to be replaceable in modules



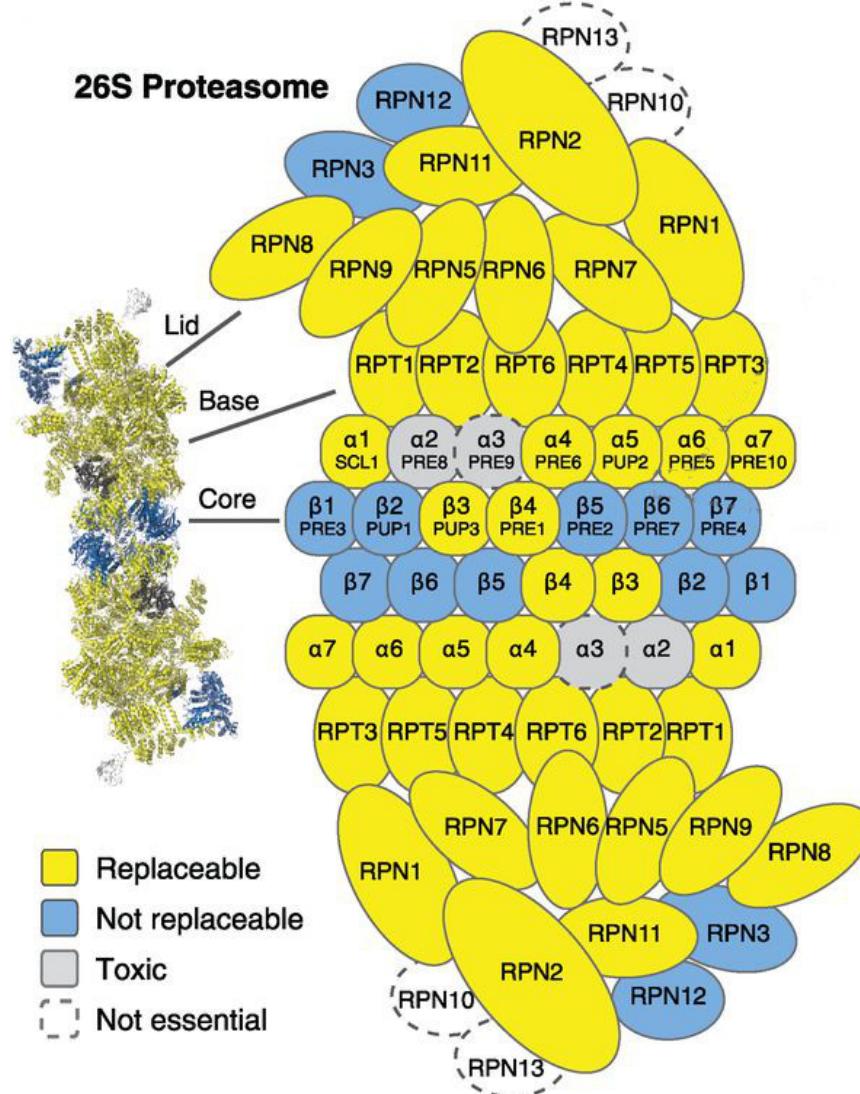
Genes tend to be replaceable in modules



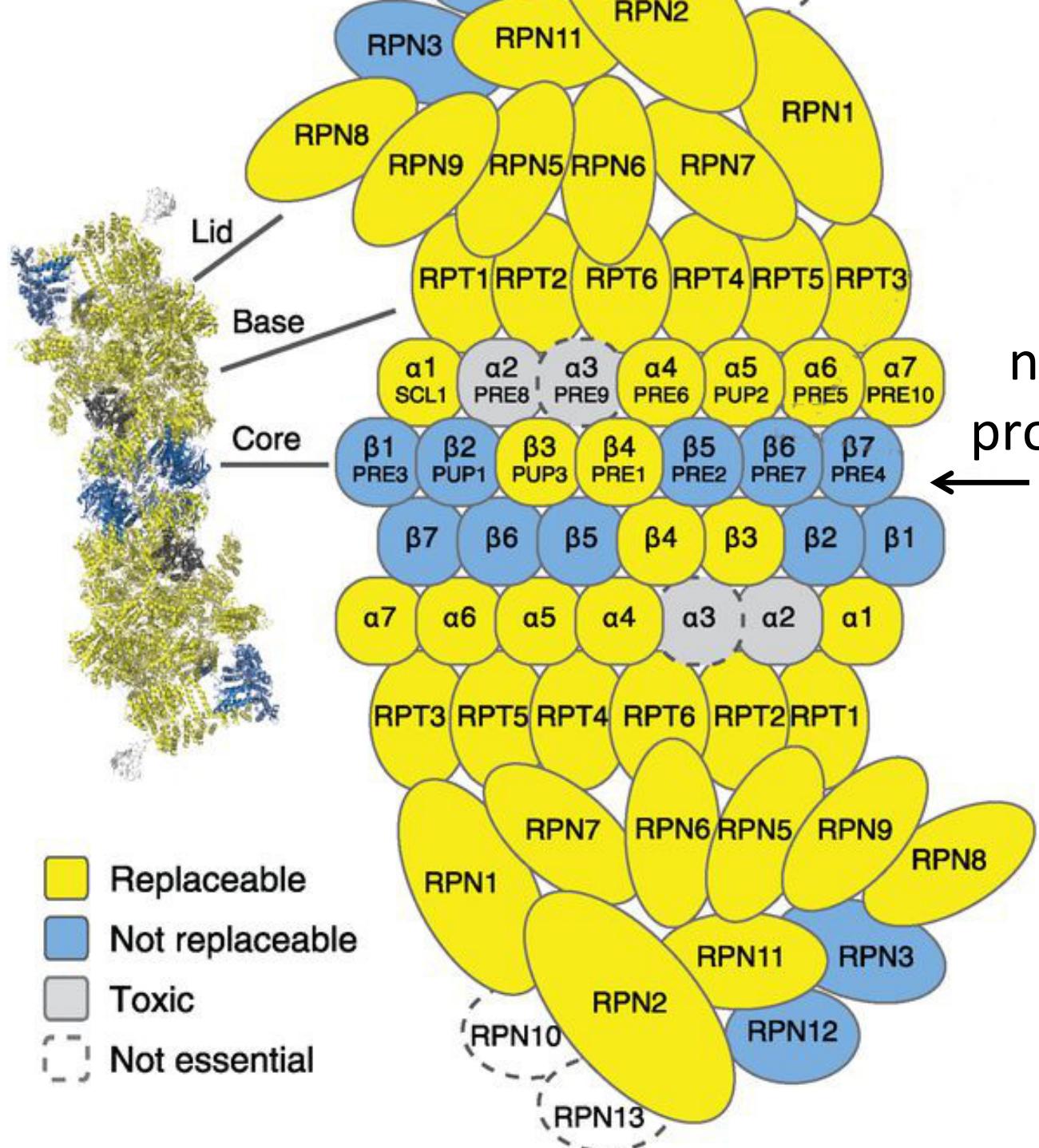
Genes tend to be replaceable in modules



Genes tend to be replaceable in modules

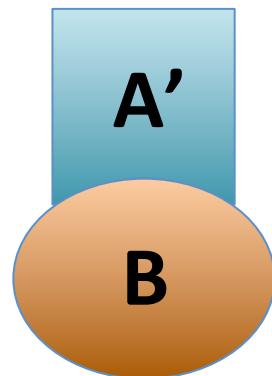
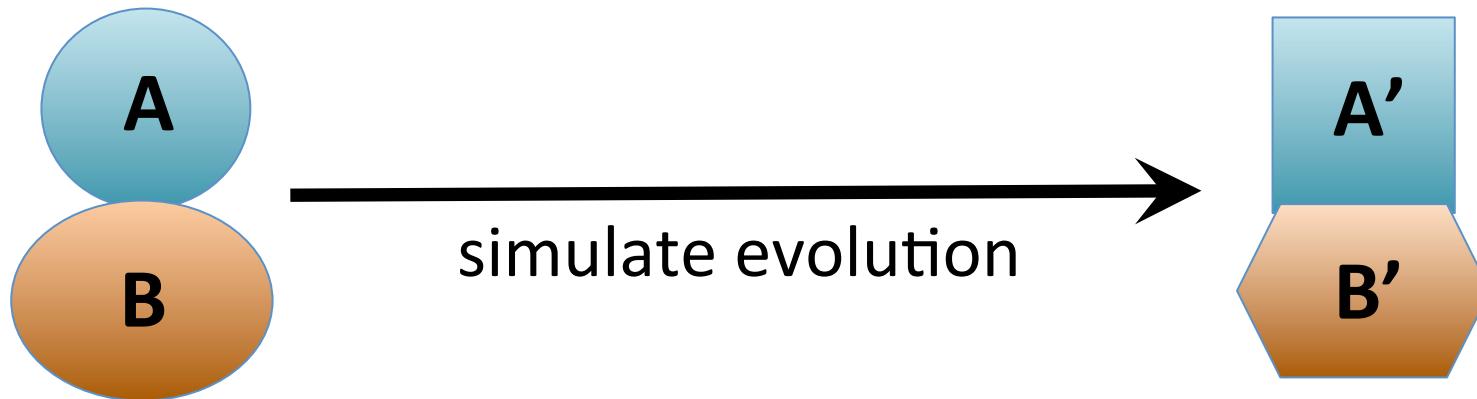


proteasome,
mostly replaceable



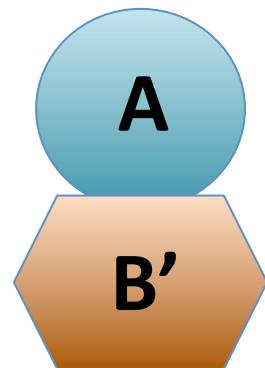
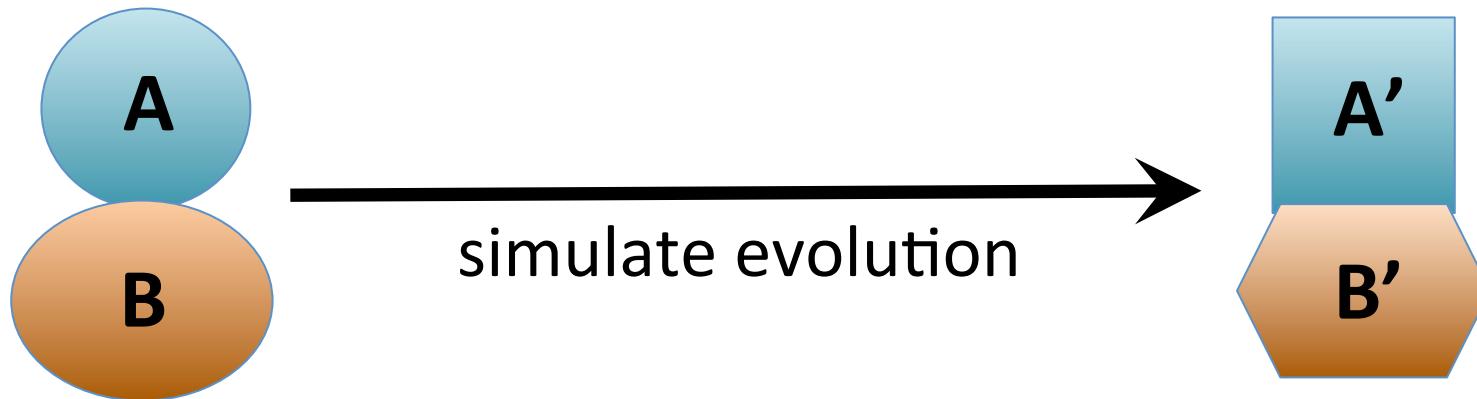
non-replaceable
proteins are mostly
in the β ring

Hypothesis: Conserved protein–protein interactions cause replaceability



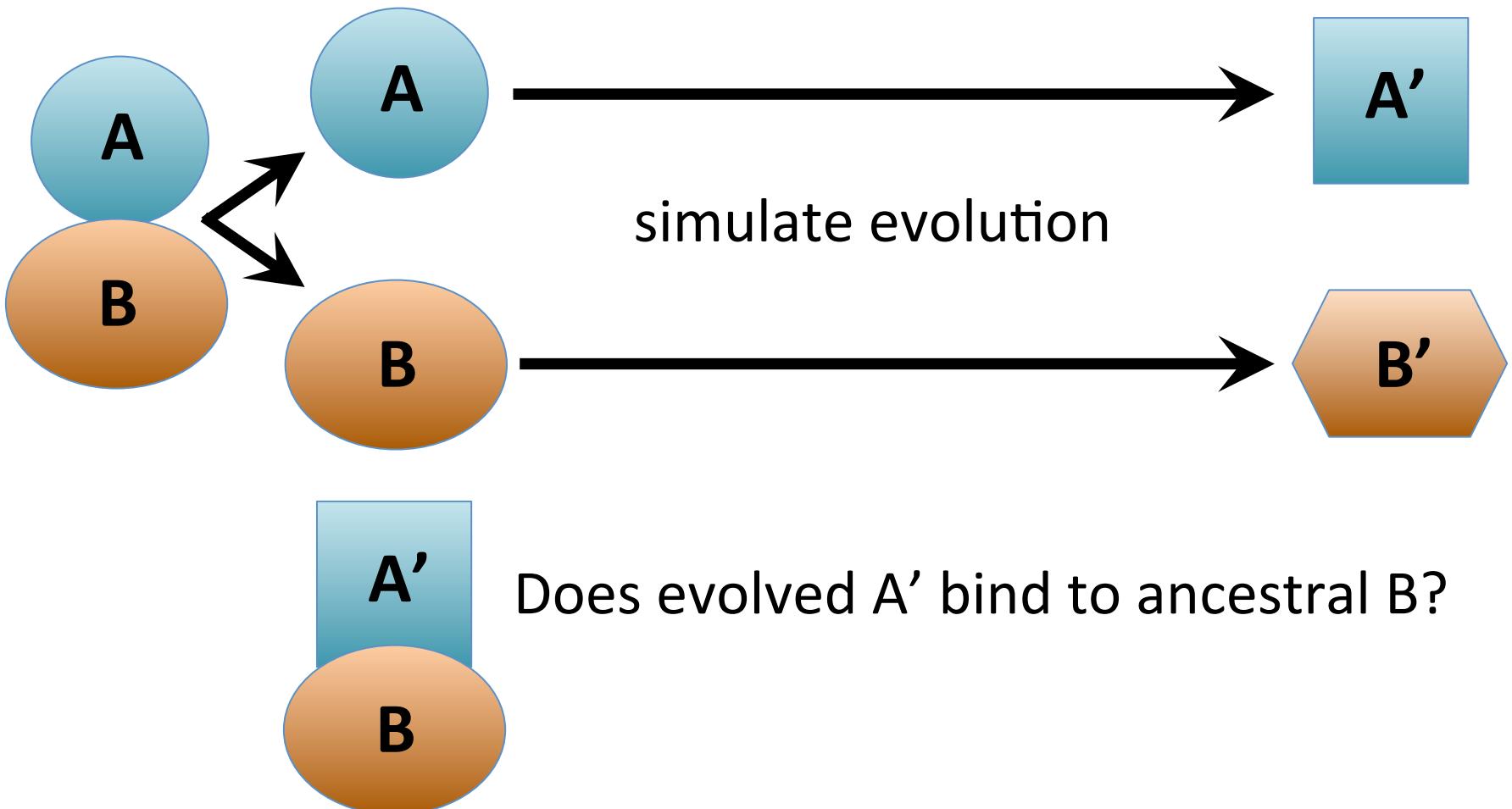
Does evolved A' bind to ancestral B?

Hypothesis: Conserved protein–protein interactions cause replaceability

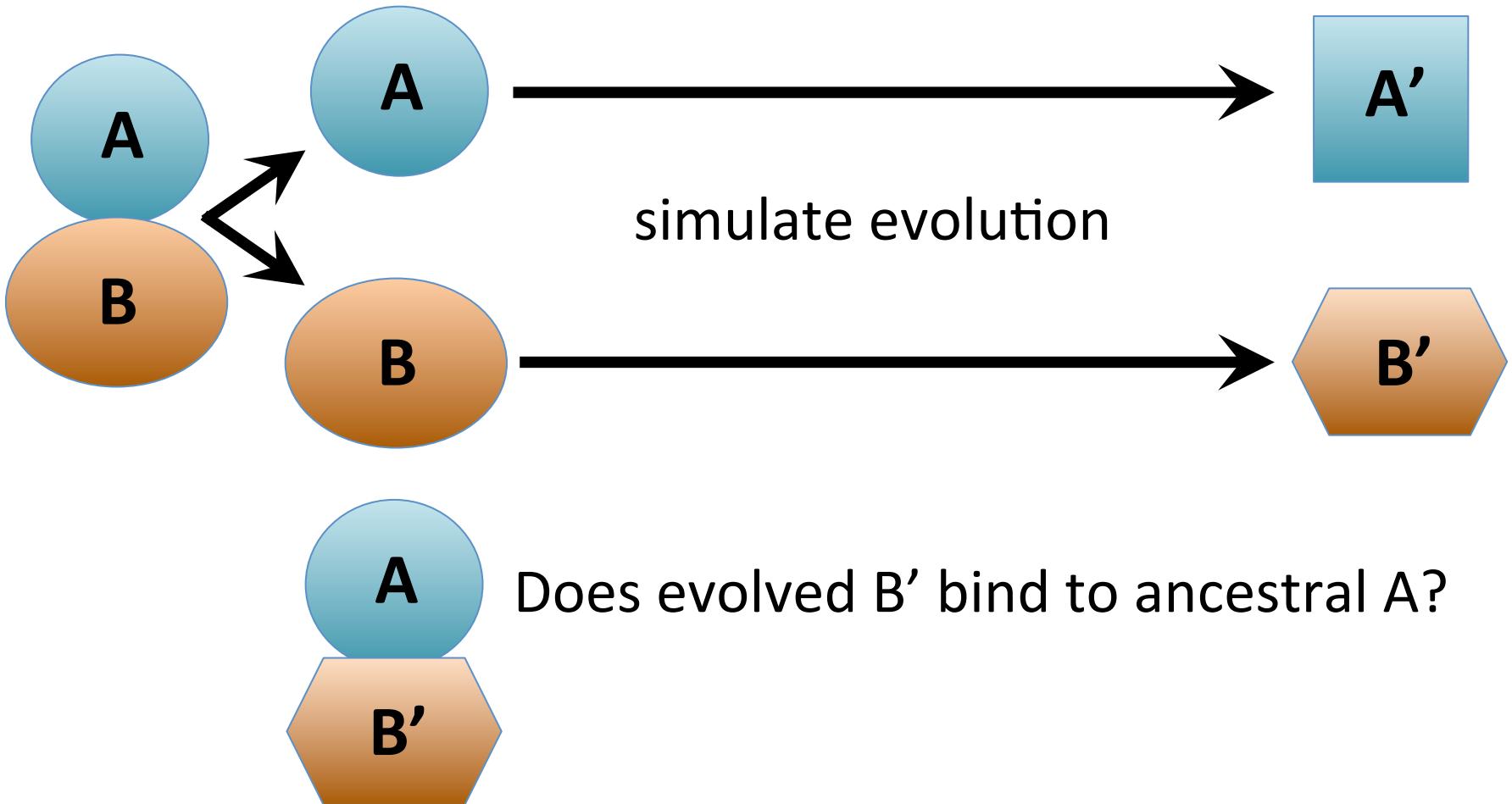


Does evolved B' bind to ancestral A ?

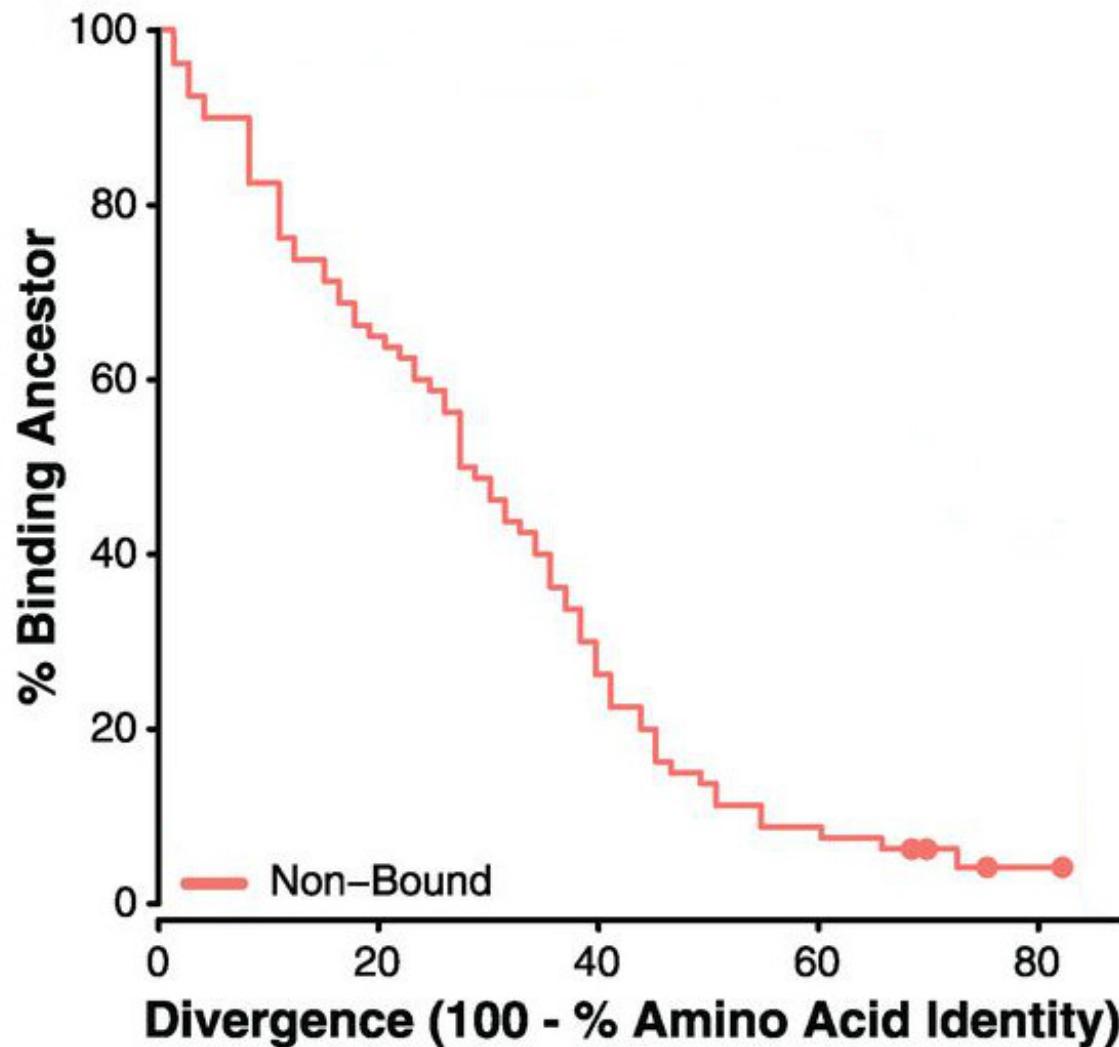
Control simulation: Evolve without binding



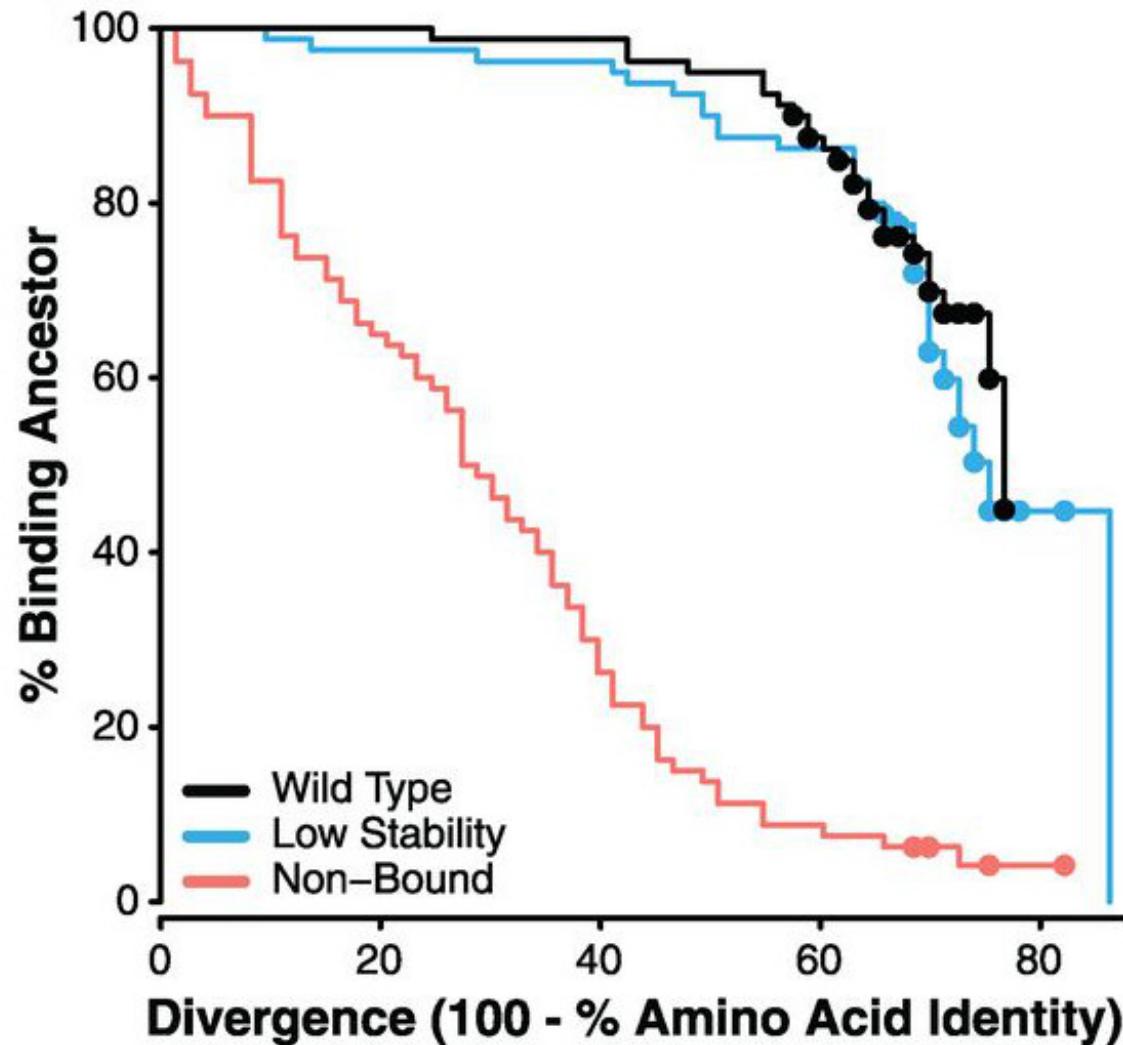
Control simulation: Evolve without binding



Binding to ancestor declines rapidly in the control scenario

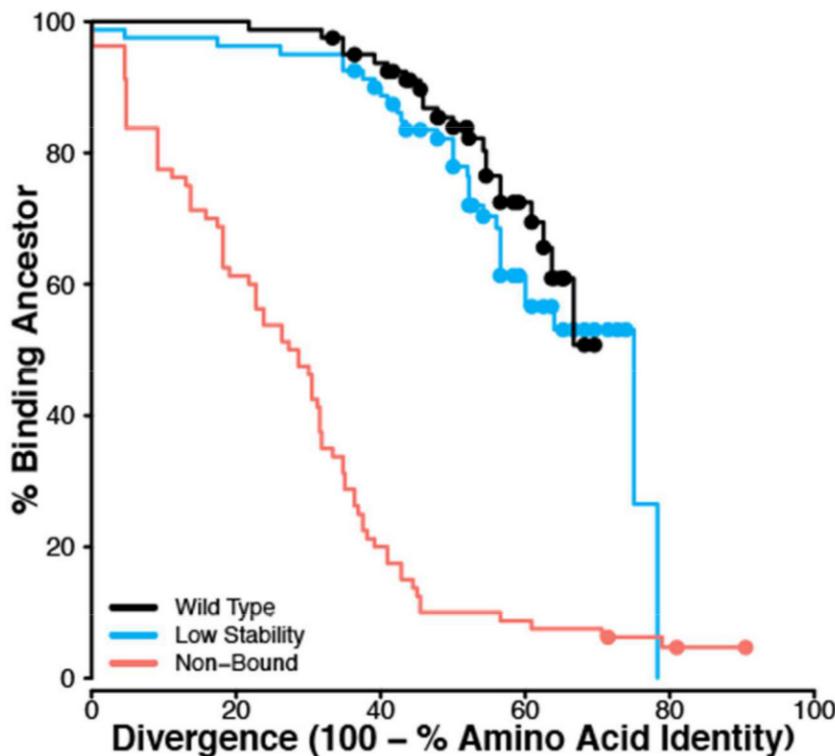


Binding to ancestor declines rapidly in the control scenario, but not in WT scenario

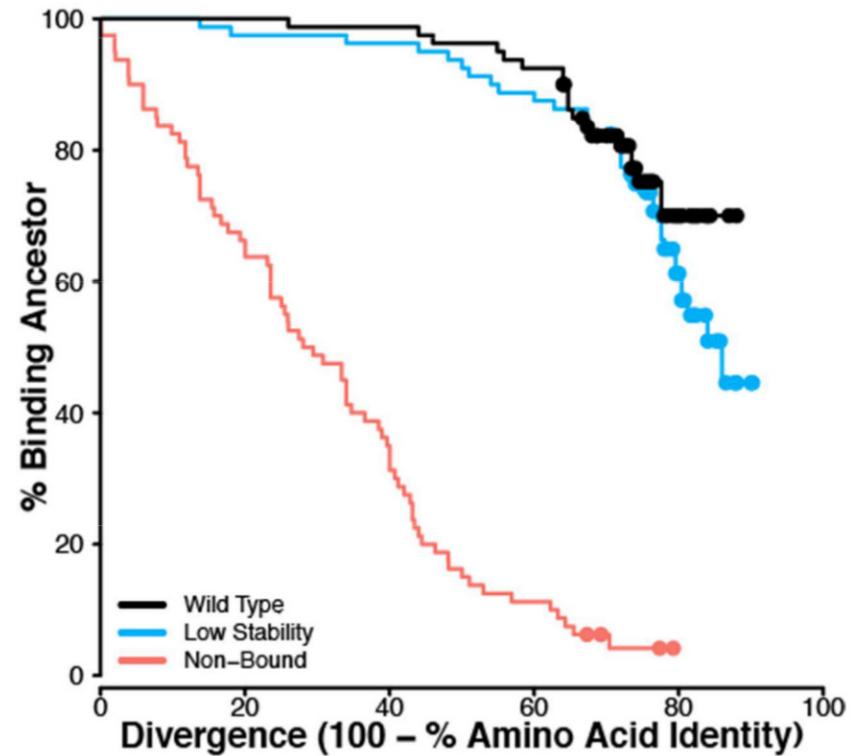


Even with substantial divergence in interface sites, evolved proteins bind ancestors

divergence in
interface sites



divergence in
non-interface sites



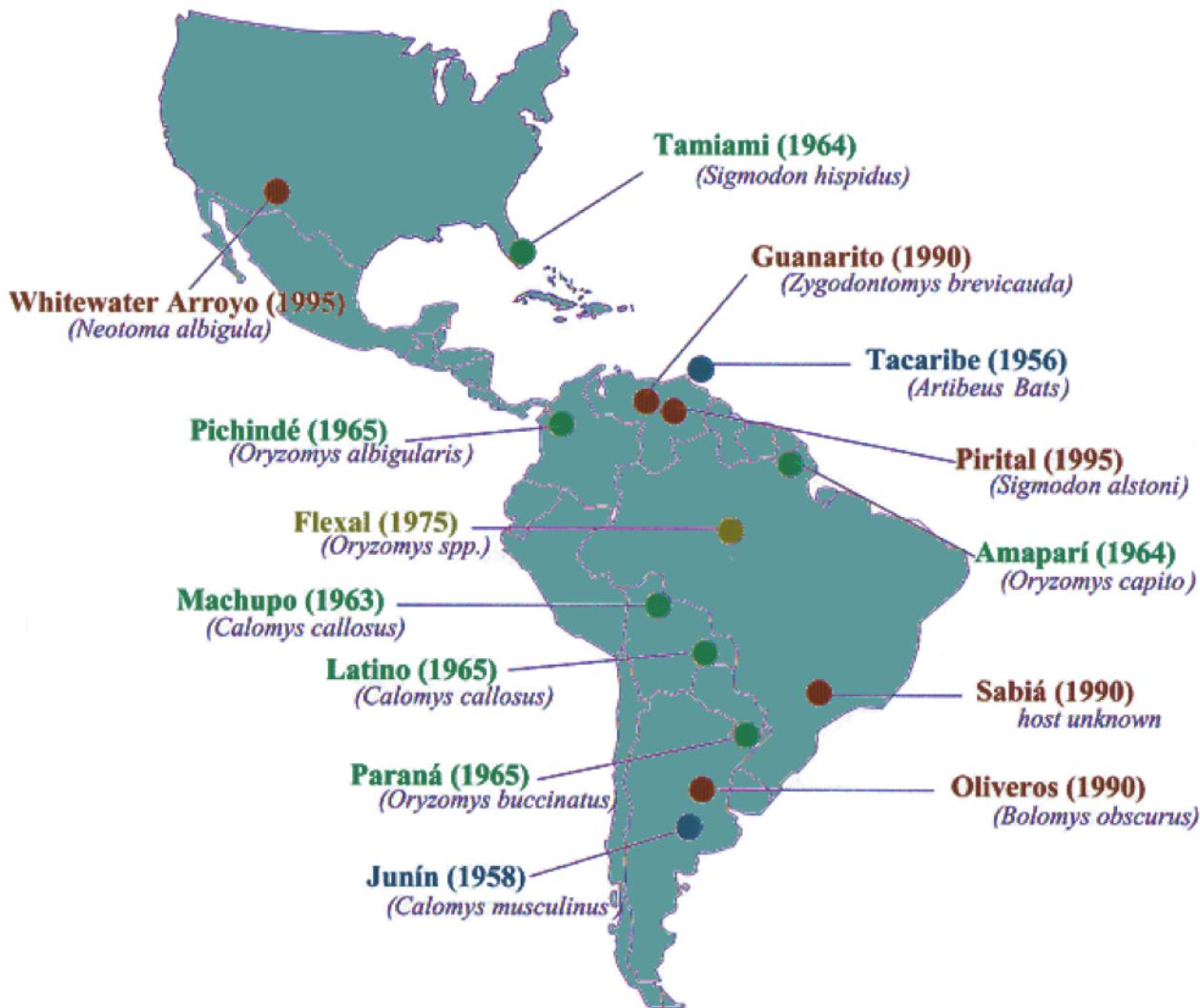
Take-home message

As protein–protein interfaces diverge, they seem to maintain biochemical similarity to their ancestral state.

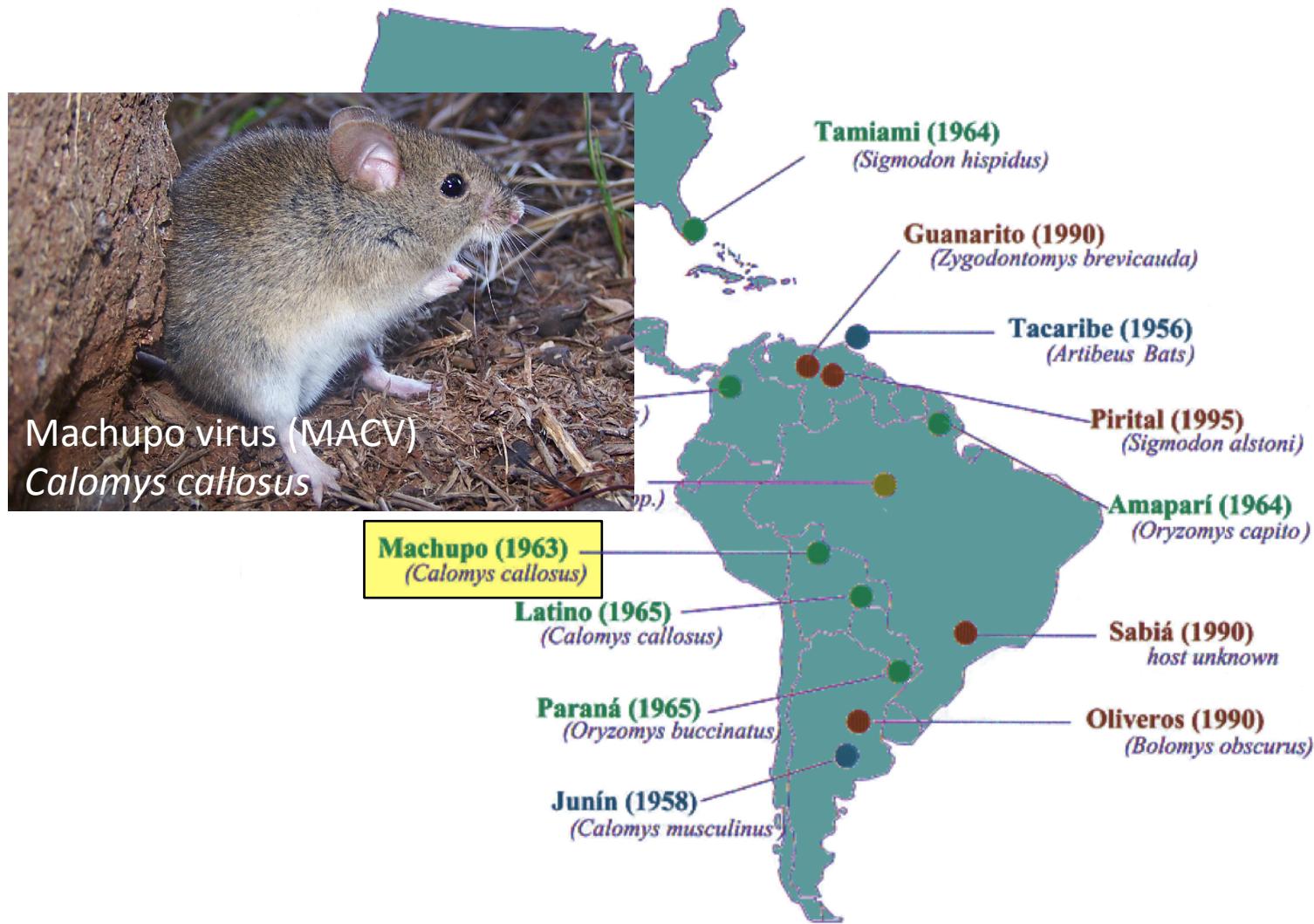
Part III: Can we predict host specificity in arenaviruses?

In collaboration with Sawyer, Ellington, and Georgiou labs
Work by Eleisha Jackson, Austin Meyer, Oana Lungu

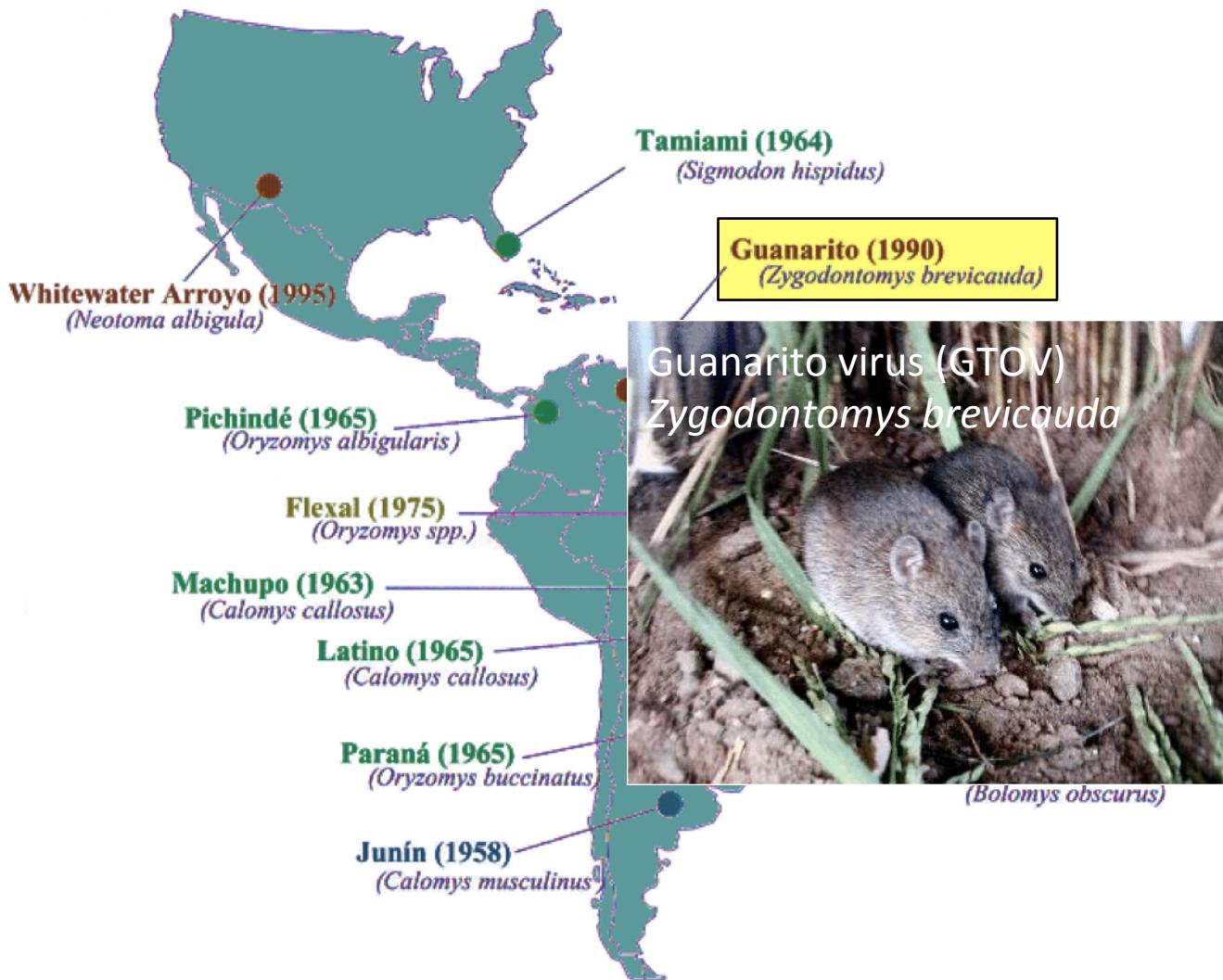
New-World arenaviruses are rodent/bat viruses found all over the Americas



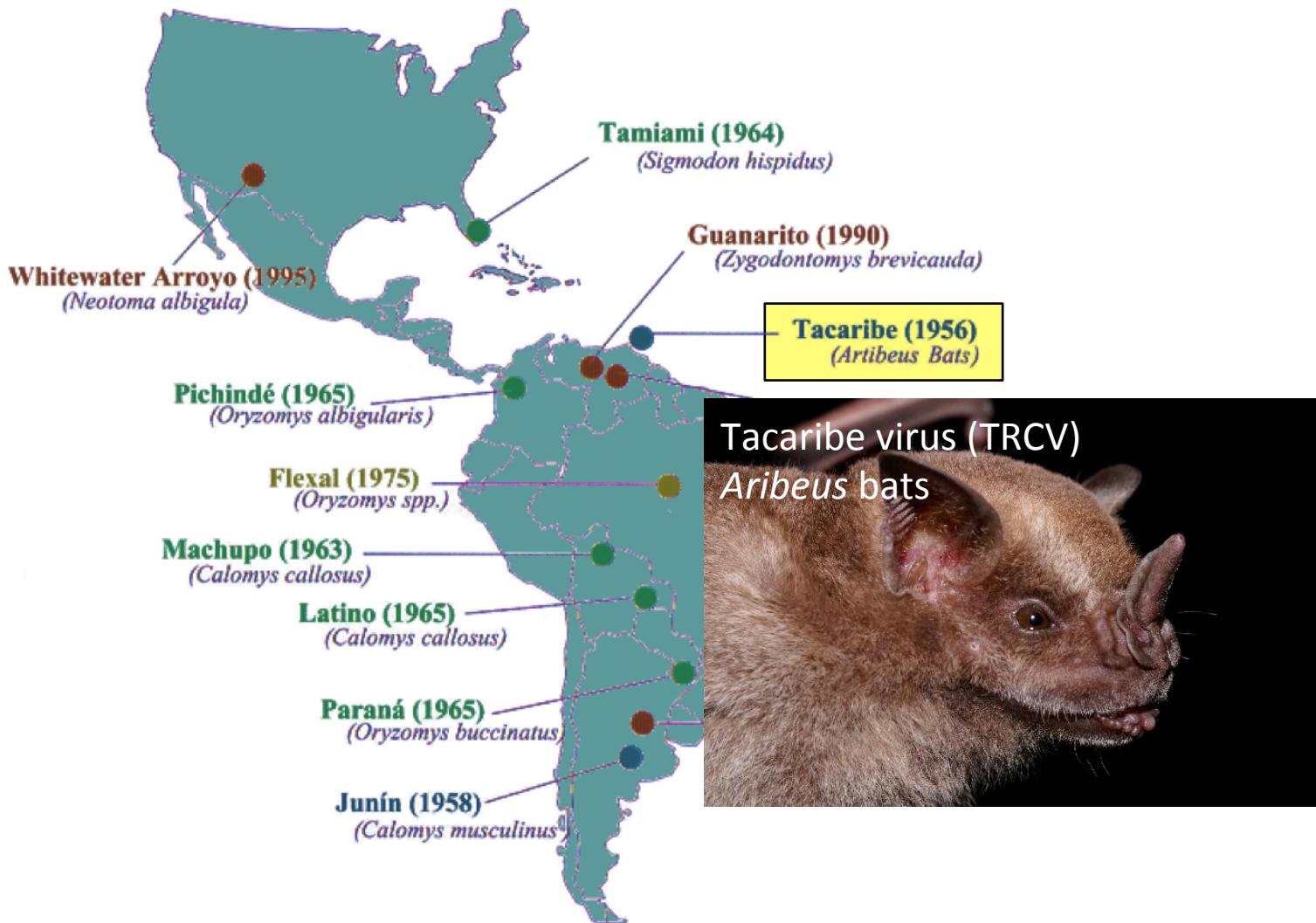
New-World arenaviruses are rodent/bat viruses found all over the Americas



New-World arenaviruses are rodent/bat viruses found all over the Americas



New-World arenaviruses are rodent/bat viruses found all over the Americas



These viruses cause severe hemorrhagic fever and high fatality in humans

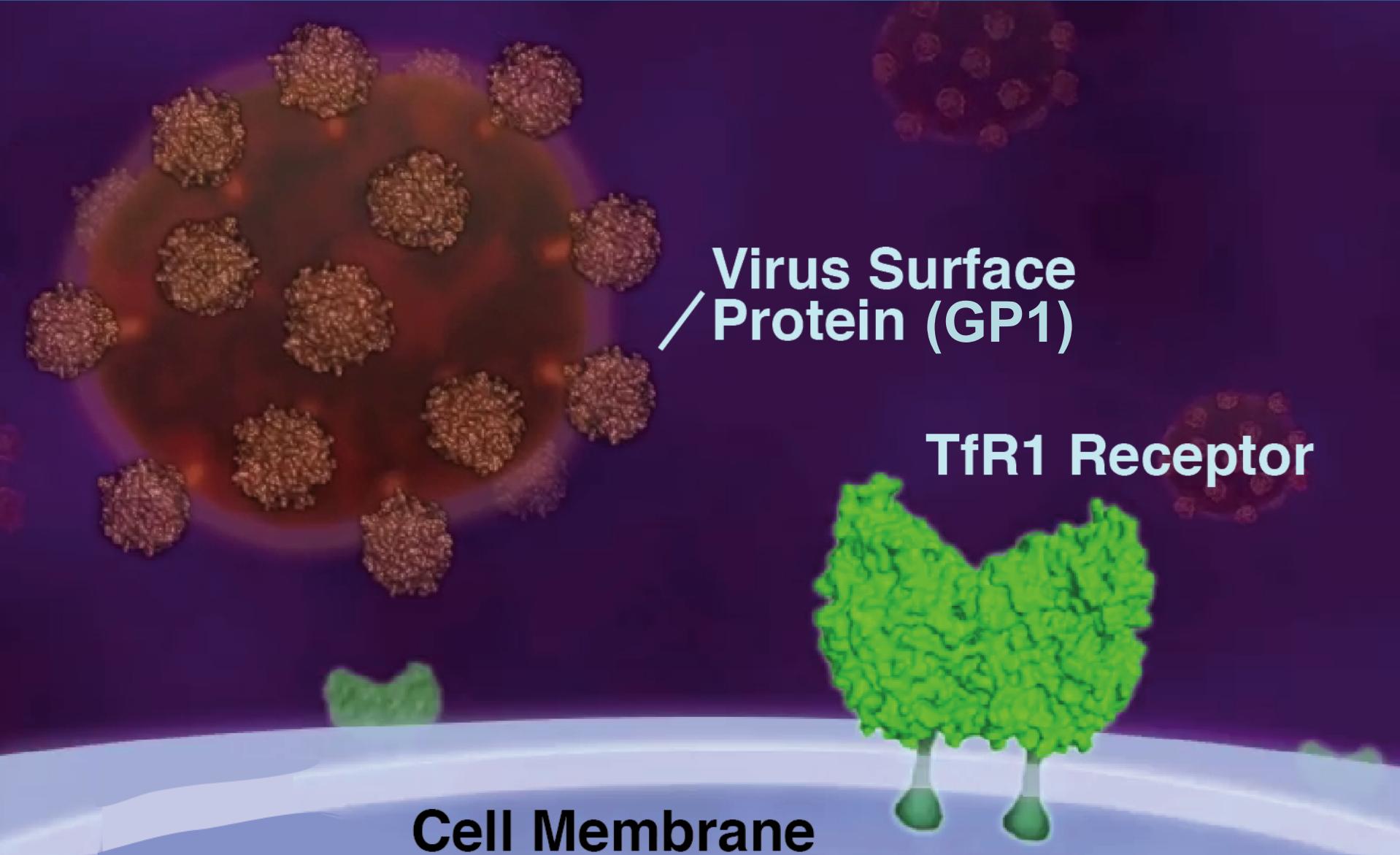


Erythematous rash in Bolivian Hemorrhagic Fever (Machupo virus)

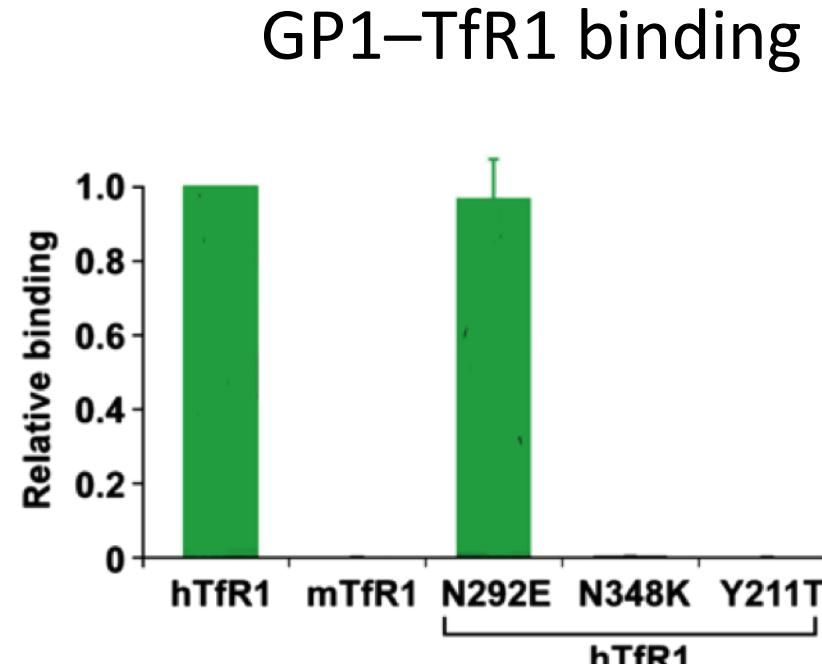
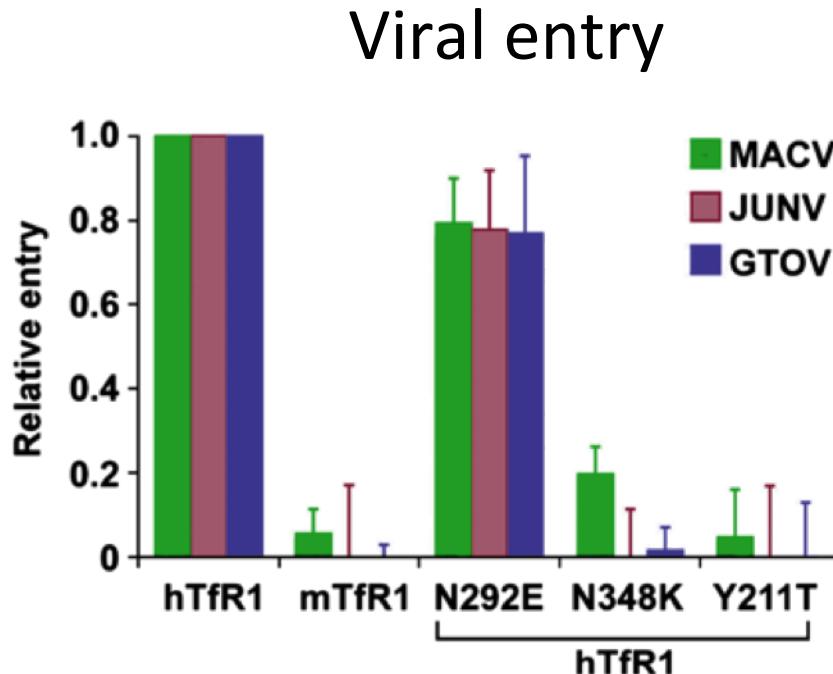


Subconjunctival hemorrhage in Bolivian Hemorrhagic Fever

Arenaviruses infect cells via the transferrin receptor (TfR1)



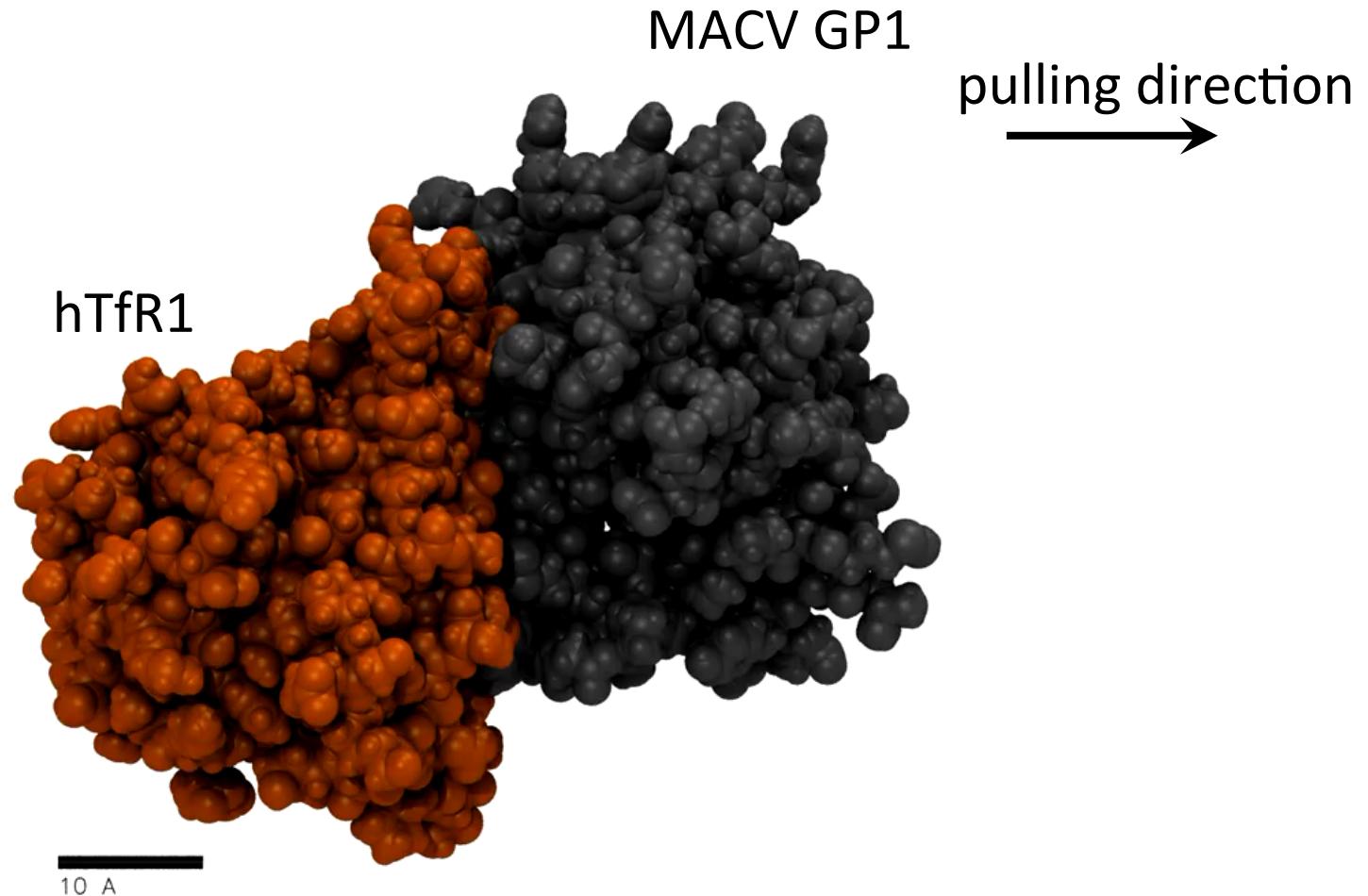
GP1 – TfR1 binding determines viral entry



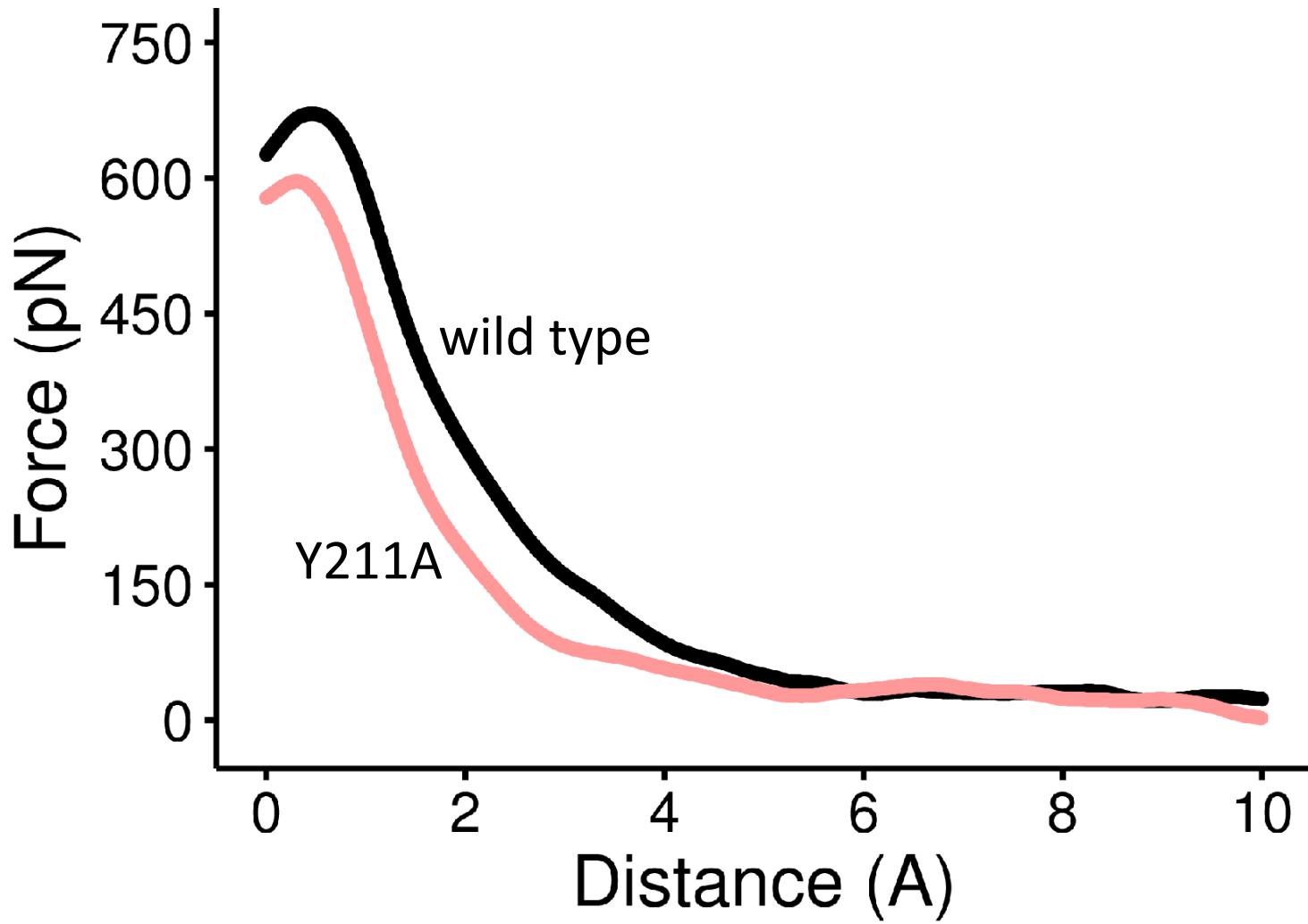
N348K and Y211T
in hTfR1 abolish binding
and entry

Can we computationally predict mutations that determine host range?

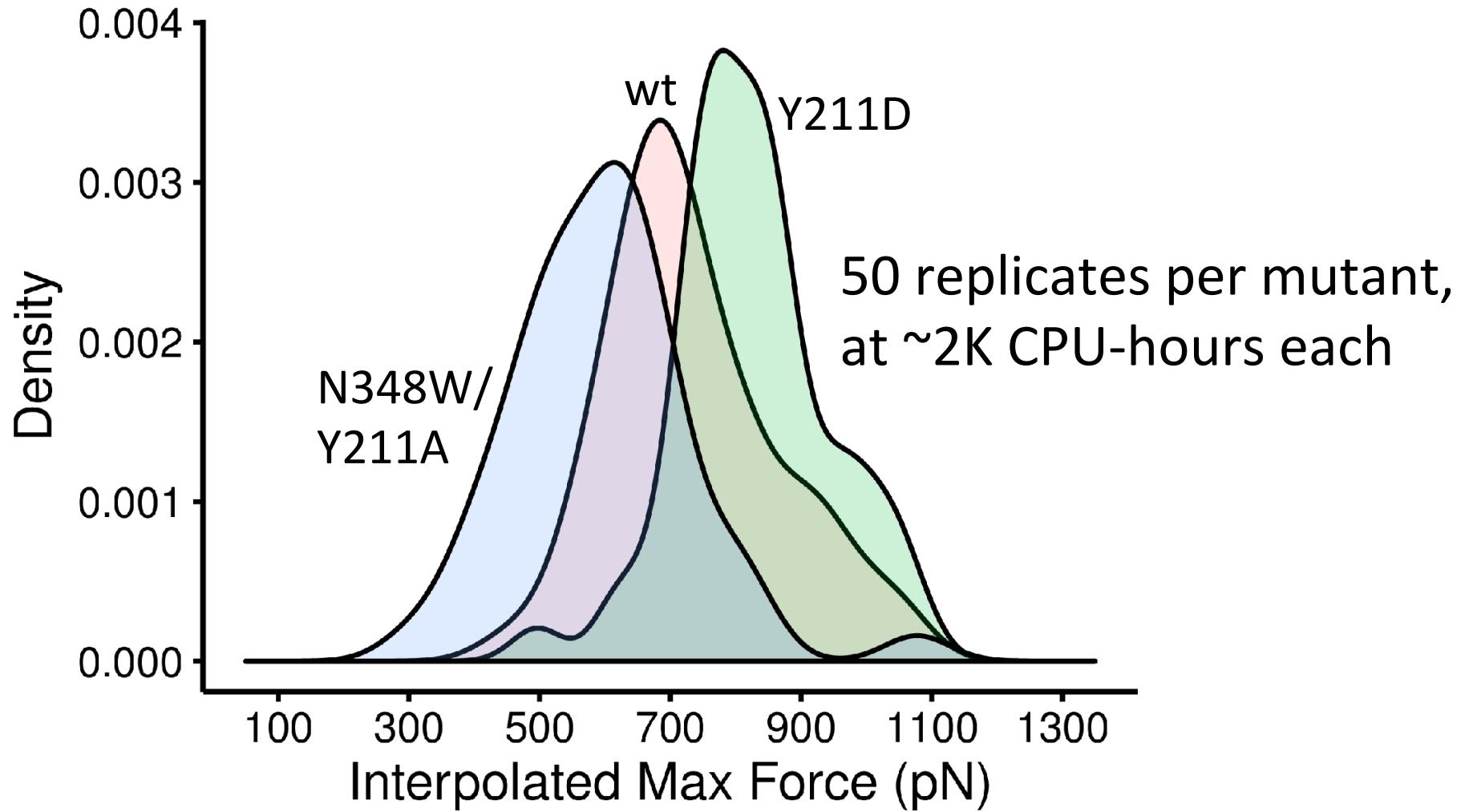
Approach 1: Steered Molecular Dynamics



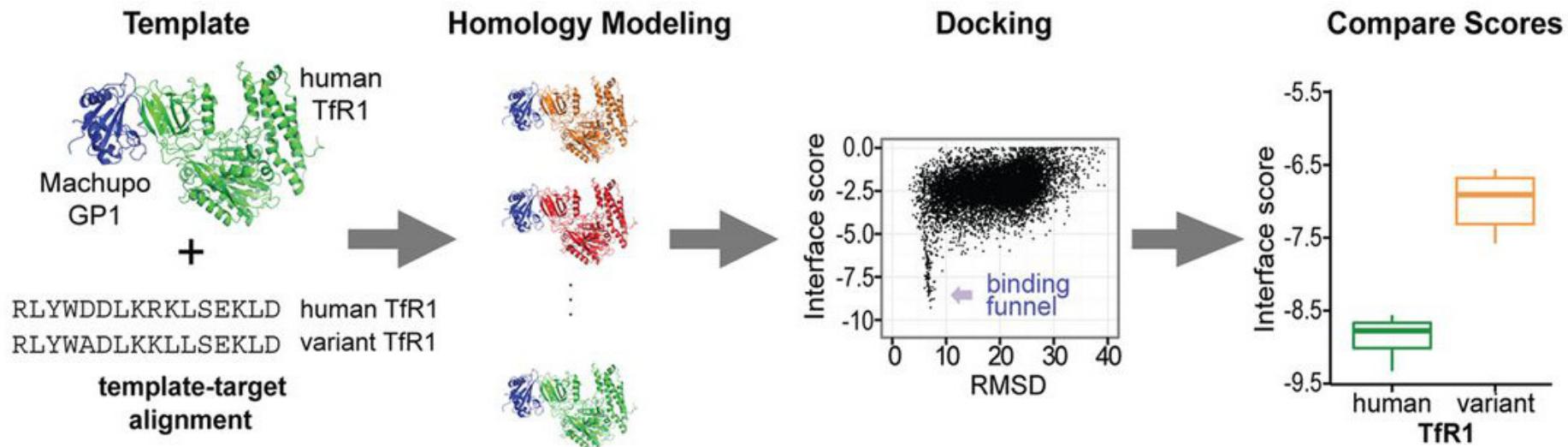
Force-distance curves differ between wild type and mutants



However, the method is noisy and requires millions of CPU-hours

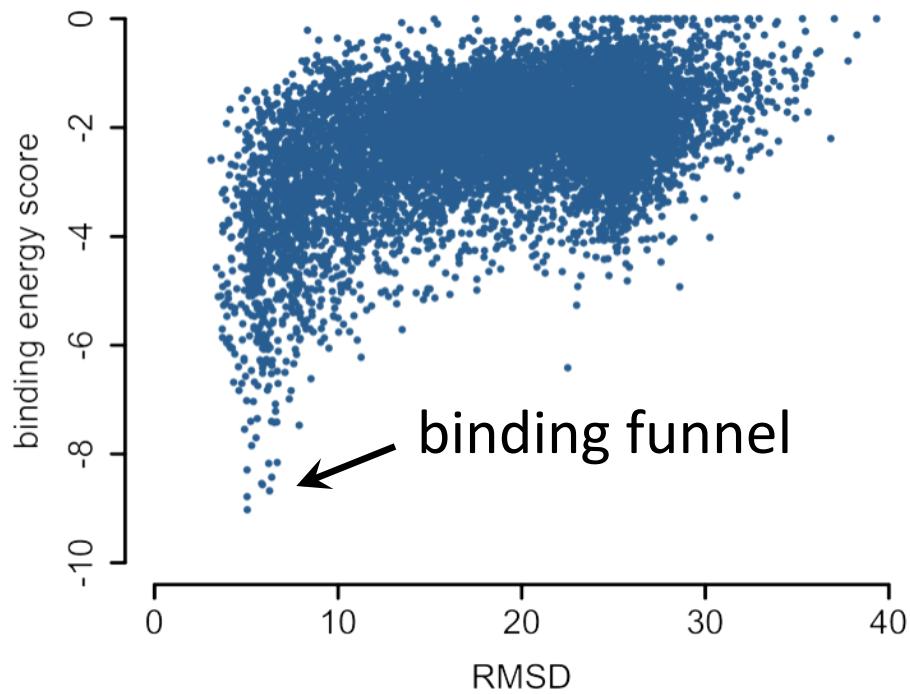


Alternative approach: Homology modeling and docking

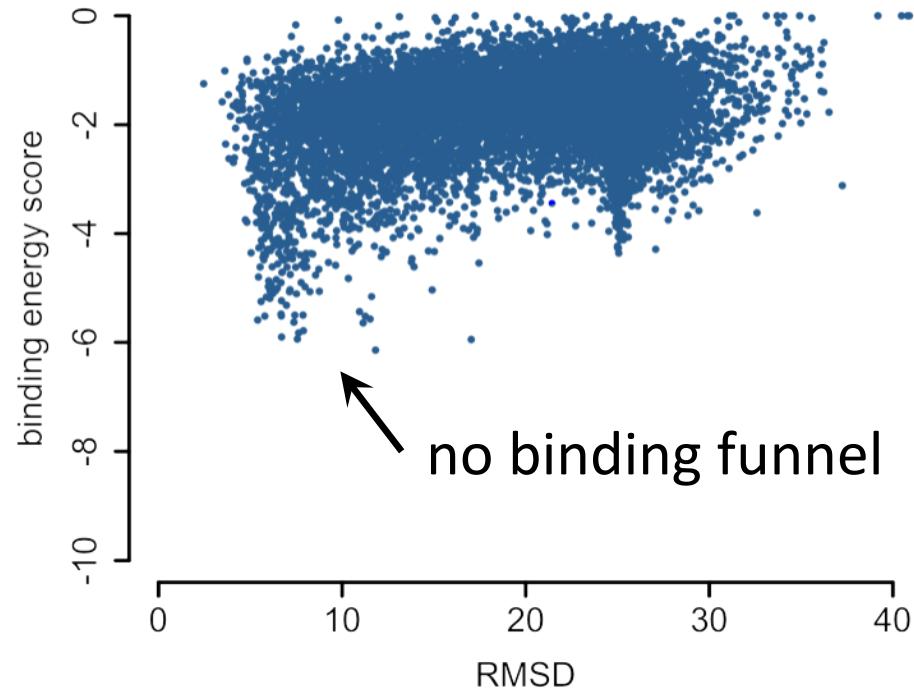


Docking results agree with experimentally observed infection patterns

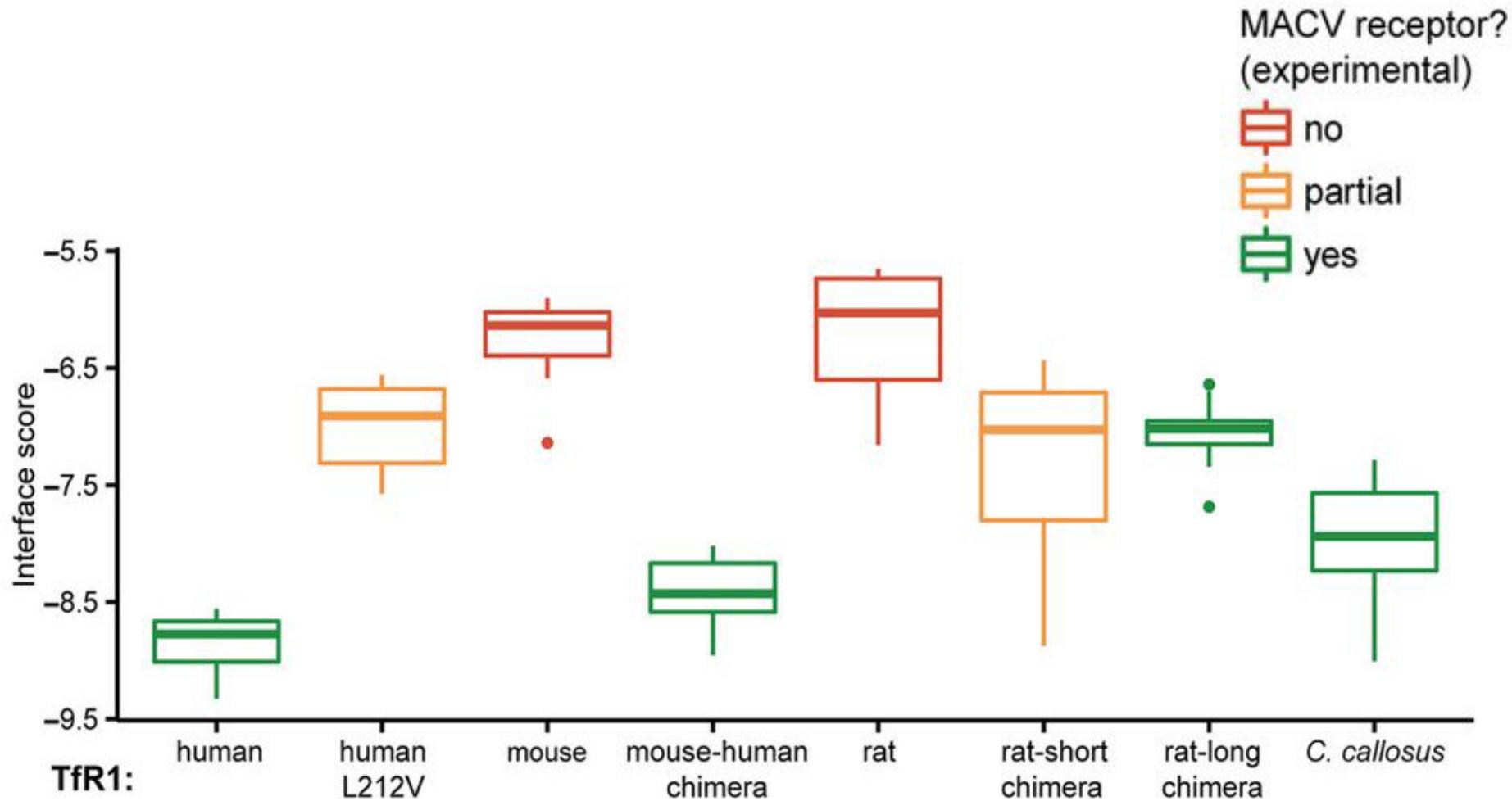
C. callosus
native host to MACV



mouse
not infected by MACV



Docking results agree with experimentally observed infection patterns



Take-home message

Computational prediction of viral host range is becoming feasible with state-of-the-art protein design.

Acknowledgments

Wilke lab

Postdocs

- Umut Caglar
- Bart Smith
- Ashley Teufel

Graduate Students

- Dakota Derryberry
- Benjamin Jack
- Eleisha Jackson
- Stephanie Spielman
- Dariya Sydykova

Former members

- Austin Meyer (now @ Texas Tech)

Collaborators

- Julian Echave
- Andy Ellington
- Aashiq Kachroo (Marcotte lab)
- Jon Laurent (Marcotte lab)
- Oana Lungu (Ellington lab)
- Edward Marcotte
- Sara Sawyer

Funding

- NSF BEACON center
- NIH
- DTRA
- ARO

