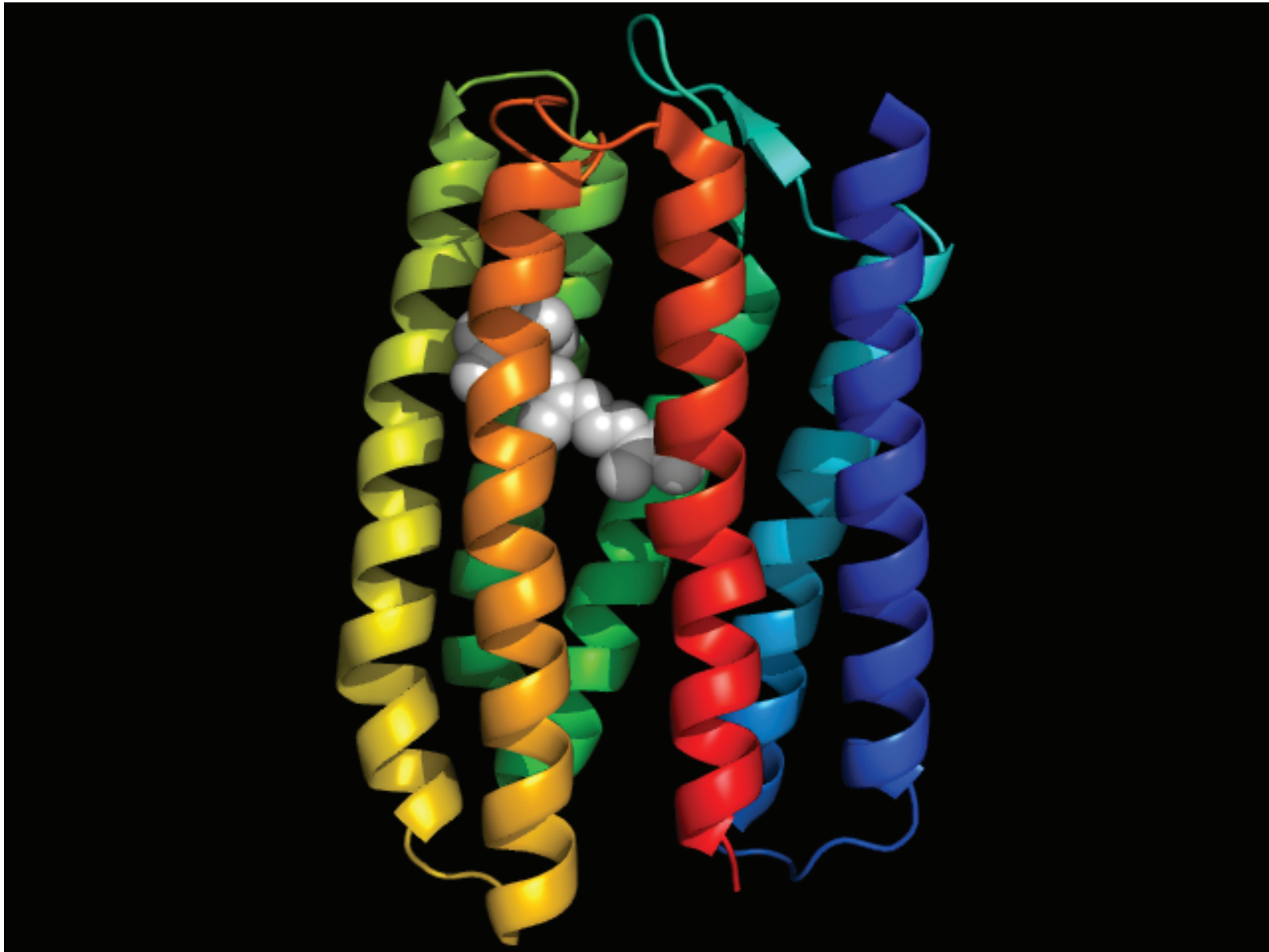


Is it worth to parameterize
sequence alignment with
an **explicit evolutionary**
model ?

Sean Eddy & E.R.



Channelrhodopsin-1

adapted from www.calvin.edu

Bacterial Rhodopsins

BACS2_HALSA
BACS2_NATPH
BACS2_HALVA
BACS1_HALSA
C7P1Y4_HALMD
D3SUL9_NATMM
BACH_NATPH
BACR_HALAR
BACR_HALSA
BACR1_HALSS
B6BSG6_9PROT
C4YP64_CANAW
B9W6Y7_CANDC
A3LUH9_PICST
B5RTR5_DEBHA
C5E3Q5_LACTC
C5DYF7_ZYGR

.TWFVWGAVGLAGTVLPI..RD
TTLFWLGAIGMLVGTLAFAWAGR
TTWFTLGLLCELLGTAVALAY.GY
ATAYLGGAVALIIVGVAFWLLYR
TTVYGLTAVVYAVALVVLWGWR
FVLLVSSIVFISAAAIFVGYSR
ASSLYINIALAGLSILLFVFMTR
AIWLWLGTAGMFLGMLYFIARGW
WIWLALGTALMGLGTLFLVKGM
TLWLGIGTLLMLIGTFYFIVKGW
GISFWVISMGMLAATAFFMETG
WAAFSVFLLLTIHLLFLYGNF
WAVFSVFALFAIVHGFYISFTDV
WALFSVFSLFAVHAFVYGFSTSS
WAVFSIFATLAVHAFVFSFTSS
WAVFSVFLVSLIYAALFVFEH
WTVTAIFGLLAVVYVLLFFVTQV

CIRHP
DAGSG
TLVPE
SLDGS
QV.SP
TLPDG
GLDDP
GETDS
GVSDP
GVTDK
NVAAG
R.KPG
R.KSG
E.KKS
R.THR
R.GTK
RNGSG

SHRRYDLVLGAGITGLAAIAYTTMG
E.RRYVTVLVTGIVSGIAAVAYVMA
ETRKRYLLIIAIPGIAIVAYALMA
PHQSALAPLAIIPVFAGLSYVGM
EHRRFCTPIVLVVALAGVASAVVA
PNQYGYAAAVA.AGSMGLAYVMA
RAKLIIVSTILVPPVSIASYTGLA
RRQKFIYATILITAIAFVNYLAMA
DAKKFYAITTLVPAIAFTMYLSML
EAREYYSITILVPGIASAAYLSMF
W.RTSVIVAGLVGTGIAFIHYMYMR
VKNSLVIPLFTNAVSVFYFTYA
LKRALLTIPLFNSAVFAFAYTYA
LKKTLVLIPLFINAVMAYTYFTYA
LKKILFIVPLFTNAIMAYCYFTYA
IHRVAVAGPLSISLVLAFSYFTMA
LSRYSLAAPFLIAFFEFFAYTYA

LGITATTVDG.....RTVY..
LGVGWVPAE.....RTVF..
LGFSGIQSEG.....HAVY..
YDIGTVIVNG.....NQIV..
AGVGTITVNG.....SEVV..
LVNGISG.....ADTD..
SGLTISVLEMPAGHFAEGSSVMLGEEVDGVVTM
LGFGLTIVEFAGEE.....HPIY..
LGYGLTMVFPGEQ.....NPIY..
FGIGLTEVQVSEM.....LDIY..
EVWVTG.....DSPT..
SNLGYAWQAVEFQH.....AGTGLRQIF..
SNLGYTWLAEFNH.....AGTGFRQIF..
SNLGTSTPTEFQH....VTTSEDLVDRQIF..
ANLGTSTRVEFNH....VSTNRLLGVRQVF..
SNLGTAVQAEFNH..LTFPNQSEVPGIRQIF..
SNLGTGTNAEFHHISVSKPVTGESPGIRQVF..

LARYIDWLVTTPL...IVLYLAML
APRYIDWILTTPL...IVYFLGLLA
VVRYVDWLLTTPL...NVWFALLA
GLRYIDWLVTTPI...LVGYVGYAA
VPLFVESMIAVGV...LYAVMARLA
LFRFLGYTAMTV...IVLVCSVA
WGRYLTWALSTPM...ILLALGLLA
WARYSDWLTTPL...LLYDLGLLA
WARYADWLTTPL...LLLDLALLV
YARYADWLTTPL...LLLDLALLA
VYRYIDWLTTPLOMVEFYLIISAVG
YAKFIWVFWGWA...VLALFEIV
YAKFVWFLGWPL...VLAIQFIV
YKWKVGYFLWPL...VLTIEVFT
YKYIGWFLWPF...VLFALIEVA
YAKYVWFLWPA...LLYLELT
YCKYIAWFLSWPI...VLFLODLA

RPG.....
GLD.....
GAS.....
GAS.....
DVE.....
GVD.....
GSN.....
GAD.....
DAD.....
KVD.....
KAN.....
TST.VLDRIENPNIFKFFLI
TNT.SFTTTEDESLLKFFISL
TQS...TDFEEDLITKFFSL
THTLESNLDAGGETVGTILSL
GVV.TRDSSNILGPRPWSFYDL
ALS...TIKRDALGSASVLDL

HMM WTVFVSGALLALVGLTLLFFVTAR RVKDG EKRRLLVILLIIPAIAAVAYVMA LGLGLTGVEAEFEH-----RQVF-- YARYIDWLLTTPL----LVLVLAELA GAD-----

HRTSAWLLAADVFVIAAGIAAAL
SREFGIVITLNTVVMLAGFAGAM
REDTVKLVVLQALTIIVFGFAGAV
RRSIIIGVMVADALMIAVGAGAVV
GRALAAIVLTPVQVIAFEVAAV
RRLTLFLFAAVLGRWLITLGSWF
ATKLFATAITFDIAMCVTGLAAAL
RNTITSLVSLDVLMICTGLVATL
QGTILALVADGIMIGTGLVGLAL
RVSIGTLVGDALMIVTGLVGLAL
SGMFWRLLLGSVVMVGGYLGEA
FQTWLVKFIIVFVYVGLLIGSI
FEALFTRVLAIEVFLGLLIGAL
FSRLFAKILATEVVFVIGLLIGAL
LSGLIVKTFATEIYVGLLIGIL
VHGLFLQICGSWFFIIGLLVGS
IHSLLVQIFGHYFWVIALVGLAL

T..T
V..P
T..P
T..D
S..G
V..D
TTSS
SPGS
T.KV
S.HT
...G
I..F
I..E
I..E
I..P
I..H
I..P

GVQ...RWLFFAVGAAGYAALYGLL
GIE....RYALFGMGAVAFGLVYVYL
SPV....SYALFAVGGALFGGVYLLY
GTL....KVALFGVSSIFHLSLFAIYL
GIV....ALIGLVVVVGGHIAIAAYLL
GTL....ALVATLGTFAALGFGLYLLF
HLM....RWFYAIACACFLVLYLILL
GVLSAGAERLVWVGISTAFLLVLLYFL
YSY....RFVWMAISTAAMLYLYVLF
PLA....RYTWLFFSTICMIVVLYFLA
YIN....ATLGFIIAGMAGVYIYEVF
STY....KFGYFTFAVFPQLLMVWVG
STY....KWGYFTFAVFPQLFAIYLV
STY....KWGYFTFSVTAQLFAEYIF
SSY....RWGYFTFAVSAQLFAMSLIL
SSY....KWGYWTMAAFAQLLVYTLIF
STY....RWGYWTIGAFMLVTEGLVL

.GTLPRALGDDPR
.GPMTESASQRSS
.RNIAVAAKSTLS
.VIFPRVVPDVPE
.GPWWTQTRGVPE
.GPVTRAAALES
.VEWAQDAAKAGT
.SSLSGRVADLPS
.FGFTSKAESMRP
.TSLRAAAKERGP
SGEAGKAAKSGN
.RDLHRSFKSPSH
.NDVVVSFGSSSH
.VNMVTAWRQSTQ
.VSMFSAAKSVHT
...KHQVLDLTI
...QRQVQALRT

VR..SLFVTLRNLTVVL...WTLYPVVWLL
GIK.SLYVRLRNLTVIL...WAIYPIWLL
DIEVSLYRFLRNFVVVL...WLVPVWVLL
QI..GLFNLLKNHIGLL...WLAYPLWLF
QRR.LLHWKARNLVFLIGMLIAVAVIALF
ERR.LLFSKLYLIVLGL...WVGL.VATGI
A...DMFNTLKLLTVVM...WLGYPVWAL
DTR.STFKTLRNLTVV...WLVPVWVLI
EVA.STFKVLRNLTVVL...WSAYPVWVLI
EVA.STFNLTALVVL...WTAYPIWLI
KALVTAFGAMRMIVTVG...WAIYPLGVF
S...NIANFFLIFFYLV...WILYPVAVGL
S...VFGNALIAFVVV...WILYPVAVGL
...KLGILLVLCQLVI...WILYPIAVGL
N...KAAIIFIAFQLLV...WILYPICWGL
S...GKLVLLVFTHVC...IYLYLVAVGL
R...GIYLILLMFMCLI...VWCYFIAWAV

SPAGIGILQ
GPPGVALLT
GAAGVGLMD
GPAGIGEAT
GVF...D
MAQGAGLAD
GVEGIAVLP
GTEGIGLVG
GSEGAGIVP
GTEGAGVVG
GYLTGGV.D
SEGNGVI.Q
SEGNGVI.Q
SEGNGKI.Q
SEGNGRI.Q
SDGNGVI.T
SEGNGKI.Q

TEMYTIVVVYLDVFSKVAFAVAVL
PTVDVALIVYLDLVTKVGFGLDAAAT
VETATIVVVYLDVTKVGFVIALLAMID
AAGVALTYVLDVLAIRVGFAGFLANLDA
AFVSLAISQYMAVLIIRVGFAGFLANLDA
DFVQQLVVIYVEVILLILGFGAIVVRSRTA
VGVTSWGYSFLDIVAKYIFAFLLNLYLTS
IGIETAGFMVIDLTAKVGFGLILLRSHGV
LNIETLLFMVLDVSAKVGFLILLRSRAI
LGIETLLFMVLDVTKVGFGLILLRSRAI
AESLNVVYNLADVFNKIAFGLVIAWAATS
PDSEAVFYGILDLLITFGLMPTILIFFAIK
PDSEAVFYGILDLLITFVGIPIILLTIAIN
PDSEAAFYGVLDFFFTFPIVGLTWLAIN
PDSEAVFYGILDLLITFSFVPIILLTWINAS
VDSSHVFFGILDLLIFVLPALLVATATATS
PDSEAVFYGILDVVFAIYPSILVWIIIV

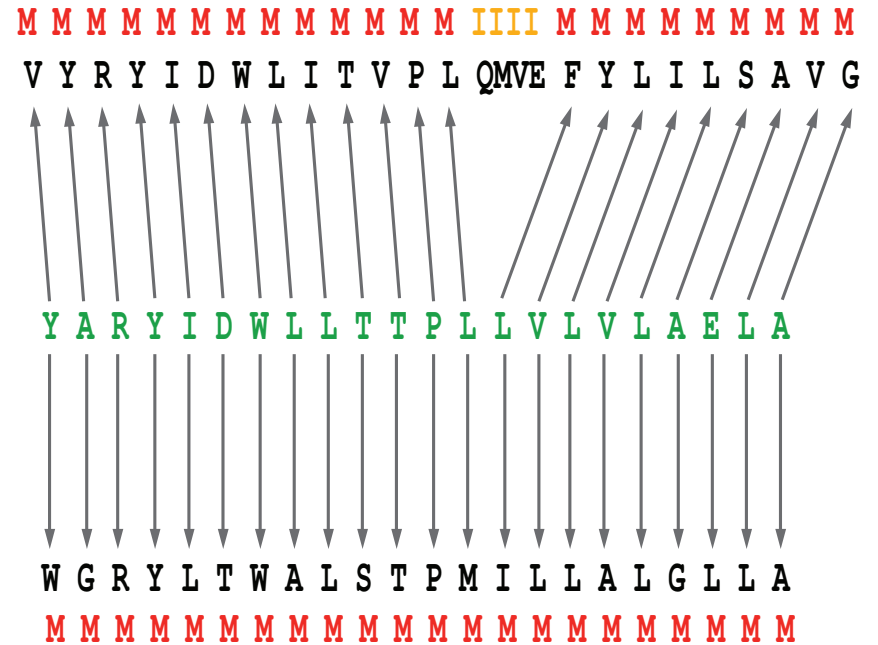
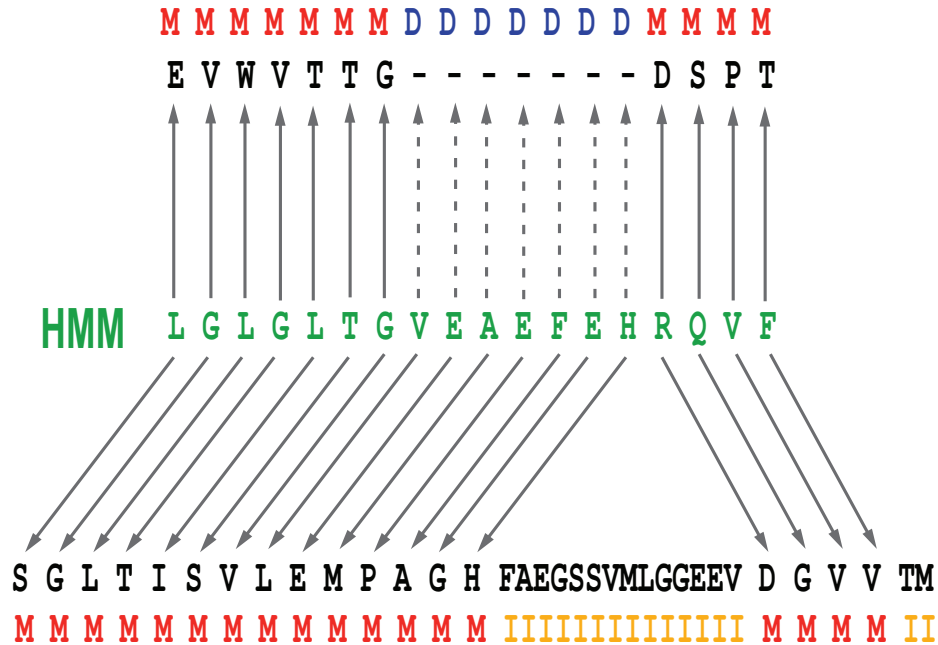
VSRLV
LRAEH
LGSAG
FMHSE
VGSAS
LSQTA
NESV
LDGAA
FGEAE
LGDTE
SSGKR
GCDEE
NVDEE
NVDEE
GVDED
SGVMP
RGEWP

RRTLLVVLADVVMVIVGVLGVAL I--E STY----RWVYFTISVAAQLVLLYLLL -GELARAAKSLSS EI--SLFNTRLRNLVVVL---WLLYPVAVLL GEEGNGI-Q ADSEAVVYGILDVLAIRVGFGLILLASATS NES--

static (fixed in time) HMM

Bacterial Rhodopsins

Alphaproteobacteria



Halobacteria (Archaea)

- **M** an evolved ancestral residue (a substitution)
- **D** a deleted ancestral residue
- **I** an insertion relative to the ancestral sequence

evolutionary (time-dependent) HMM

Homology is an evolutionary question

Homology detection is hypothesis testing

Forward score

Posterior of \mathbf{H} given s

$$F = \log \frac{P(s | \mathbf{H})}{P(s | \mathbf{R})}$$

$$P(\mathbf{H} | s) = \frac{e^{F+\rho}}{1 + e^{F+\rho}}$$

Evolutionary distance is a nuisance parameter in

$$P(s | \mathbf{H})$$

Current approaches assume
(implicitly or explicitly)

a fixed evolutionary distance

An explicit time-parameterization
allows to

Integrate over Evolutionary Distance

Homology detection

Homology coverage

Optimize for Evolutionary Distance

Alignment of homologs

Affine Gap cost

A way of dealing with variability



$$S(A,A) + S(V,H) + S(G,G) + \beta + \eta + \eta + S(V,V) + S(L,L) + \beta + S(K,K)$$

substitution
matrix

gap open cost	β
gap extent cost	η

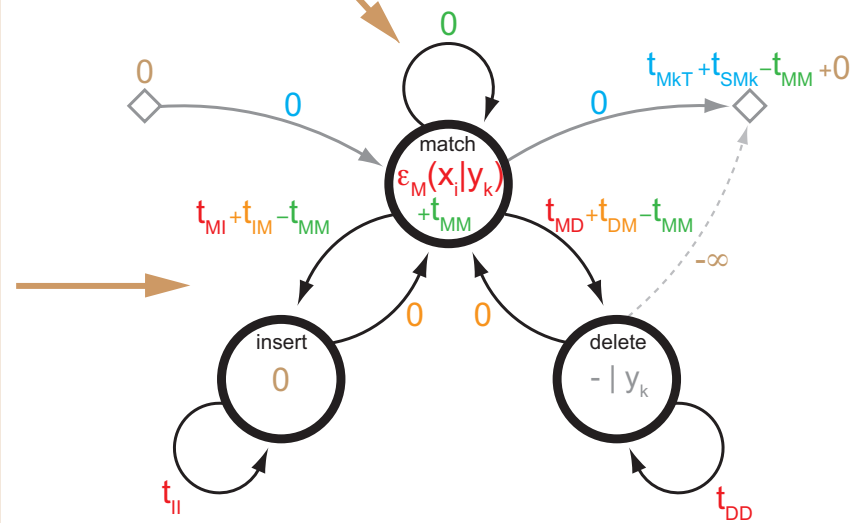
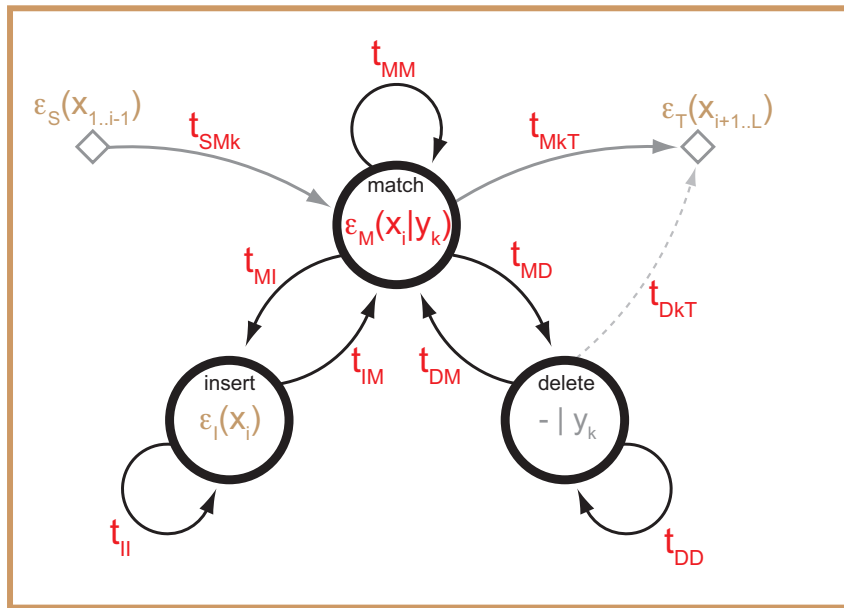
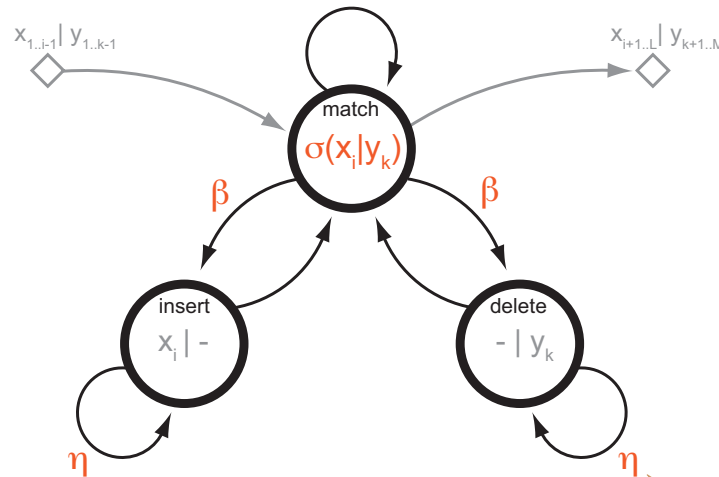
affine gap cost

BLAST syncs (empirically) the choice of substitution matrix with that of the affine gap costs

substitution matrix	BLOSUM62
gap open	-11
gap extent	-1

HMMs formalize sw-like affine methods

From Smith-Waterman to an HMM



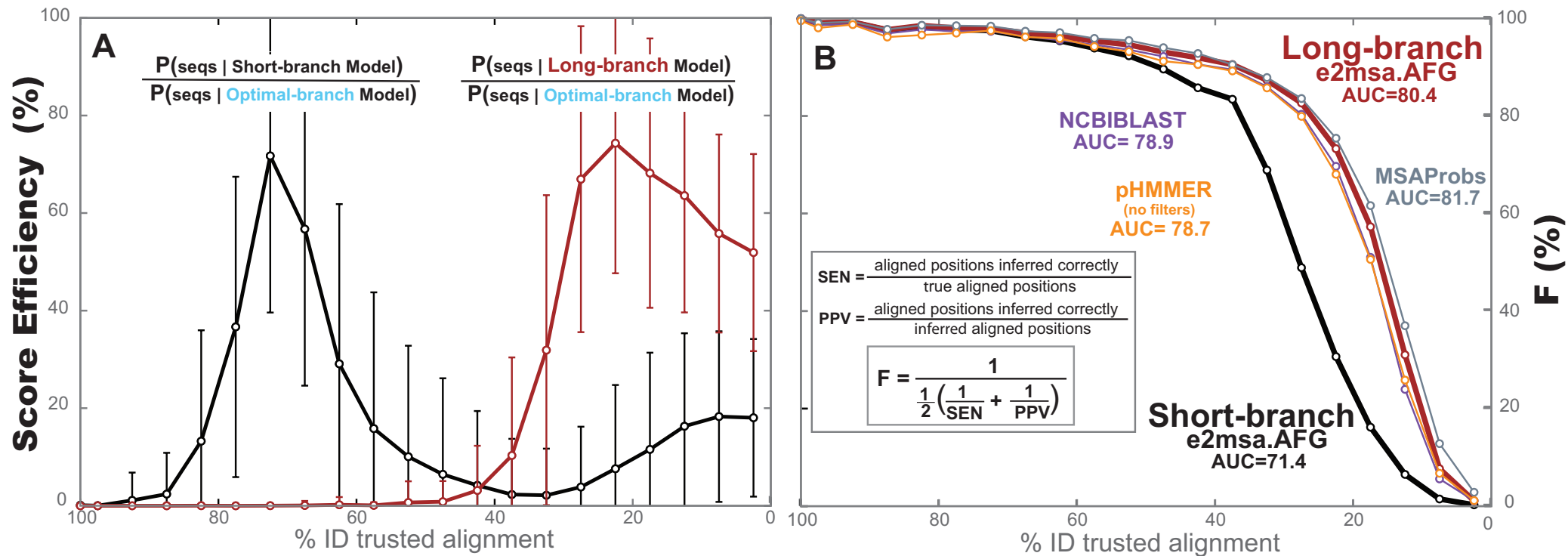
Eddy & Castellano *unpublished*

A probabilistic evolutionary model provides time-dependent HMM transitions

Is it worth to
parameterize pair and
profile HMMs with an
explicit evolutionary
model?

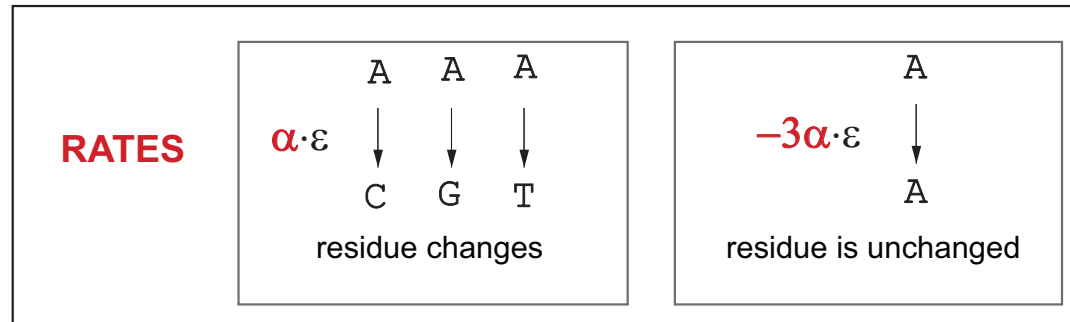
Alignment Accuracy Benchmark

Global Homology Set



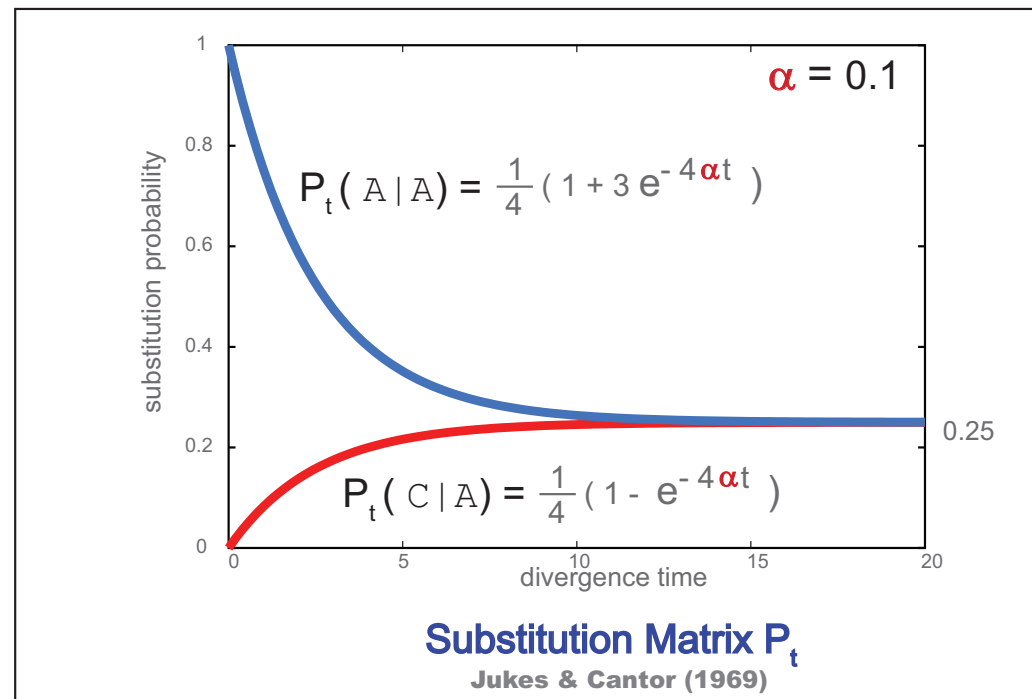
Evolution of residue substitutions

Assume
For very small times ε :



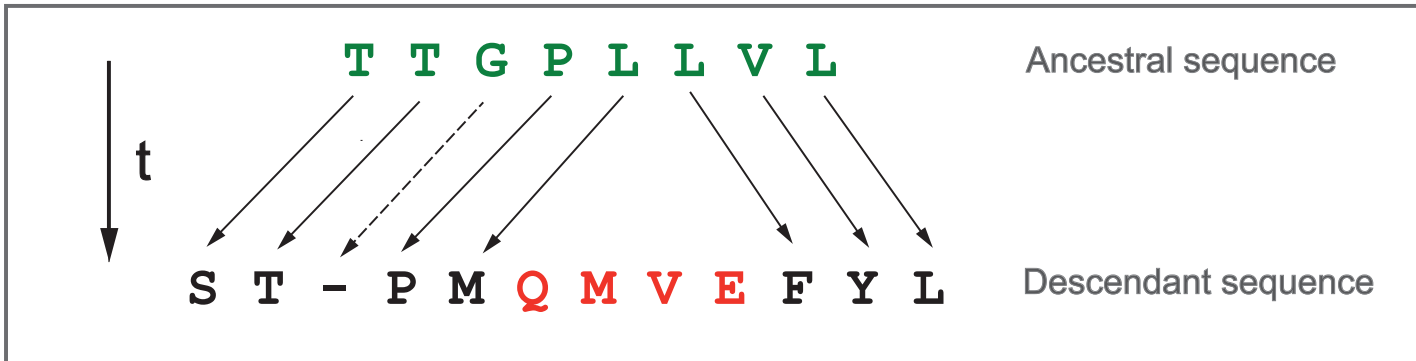
propose and solve differential equations

Infer
For finite time t :



Evolution of Insertions

compatible with affine models



Substitutions

infinitesimal rate α



$$P_t(S | T)$$

Insertions & Deletions

rate for deleting an ancestral residue μ_A

rate for starting a new insert with "n" residue $\lambda (1 - s_I) s_I^{n-1}$

rate for deleting a whole insert with "n" residues $\mu (1 - s_D) s_D^{n-1}$

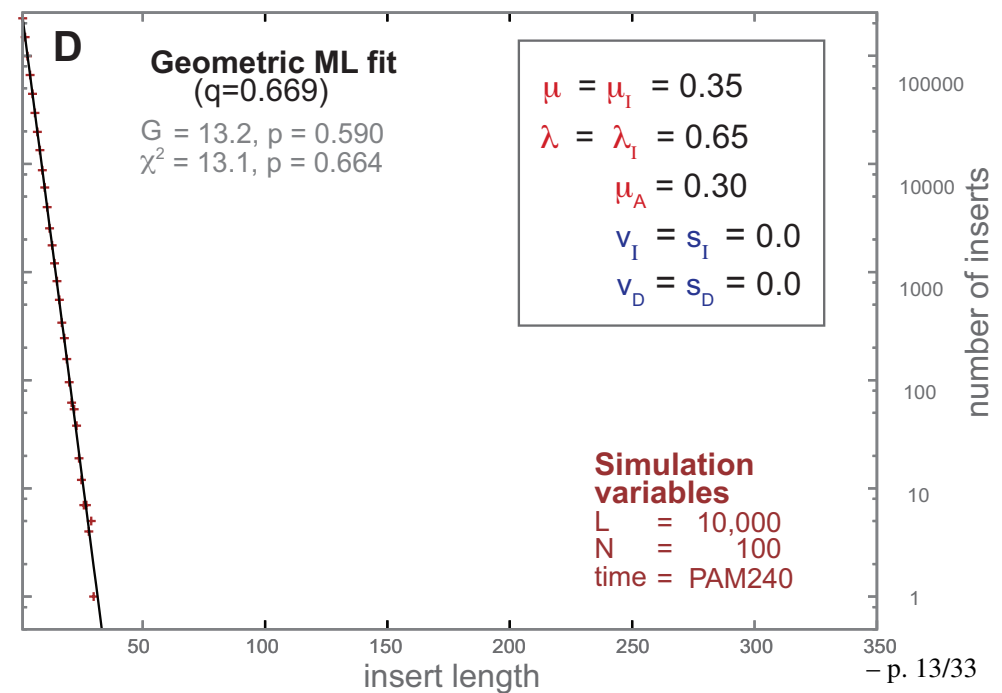
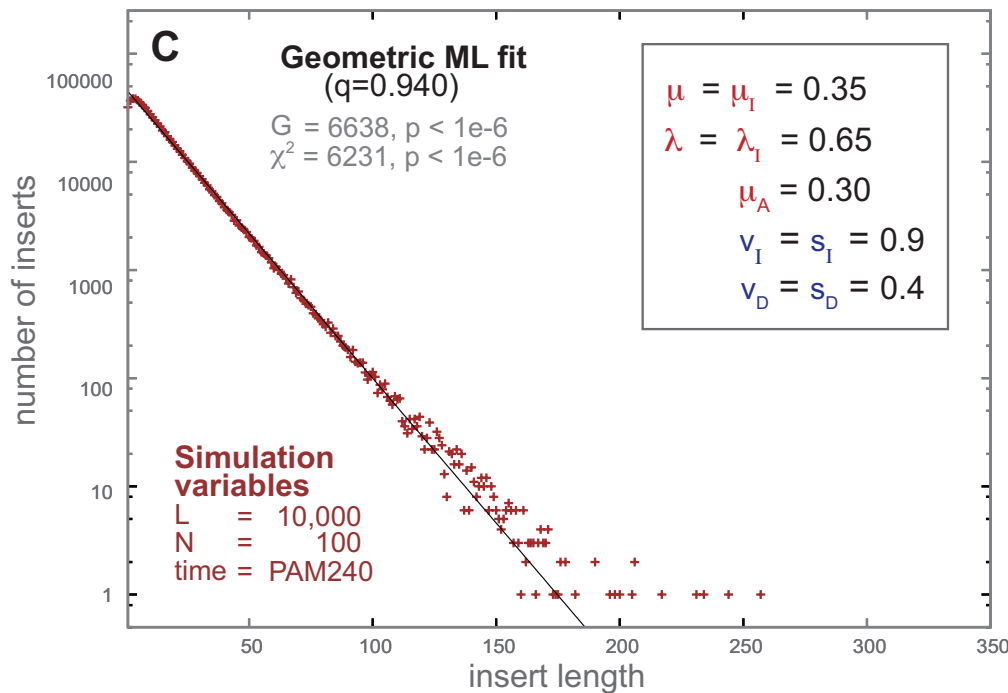
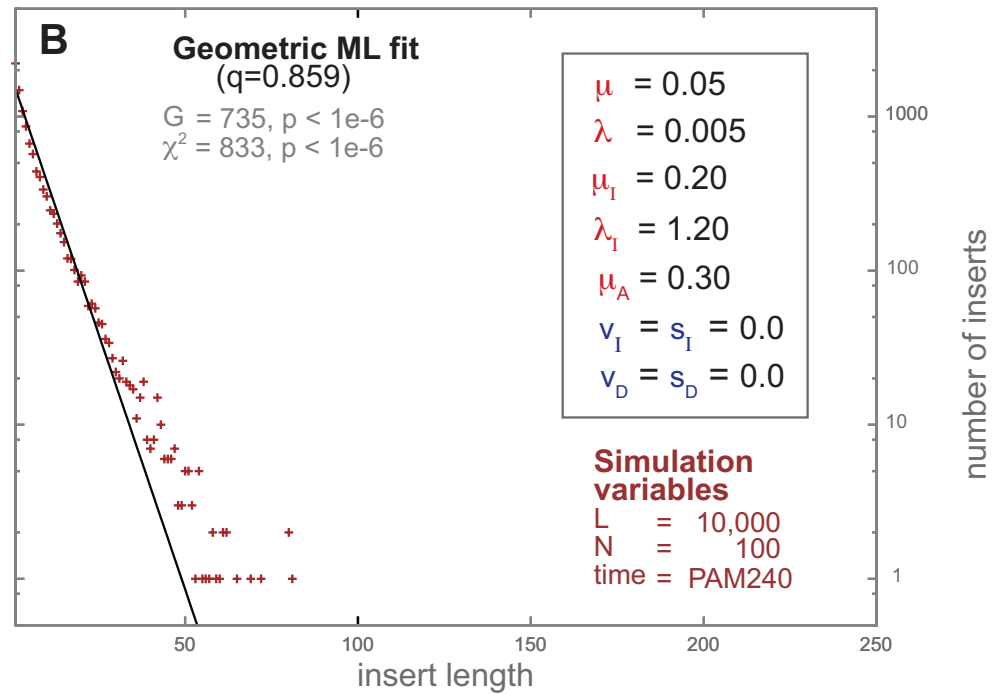
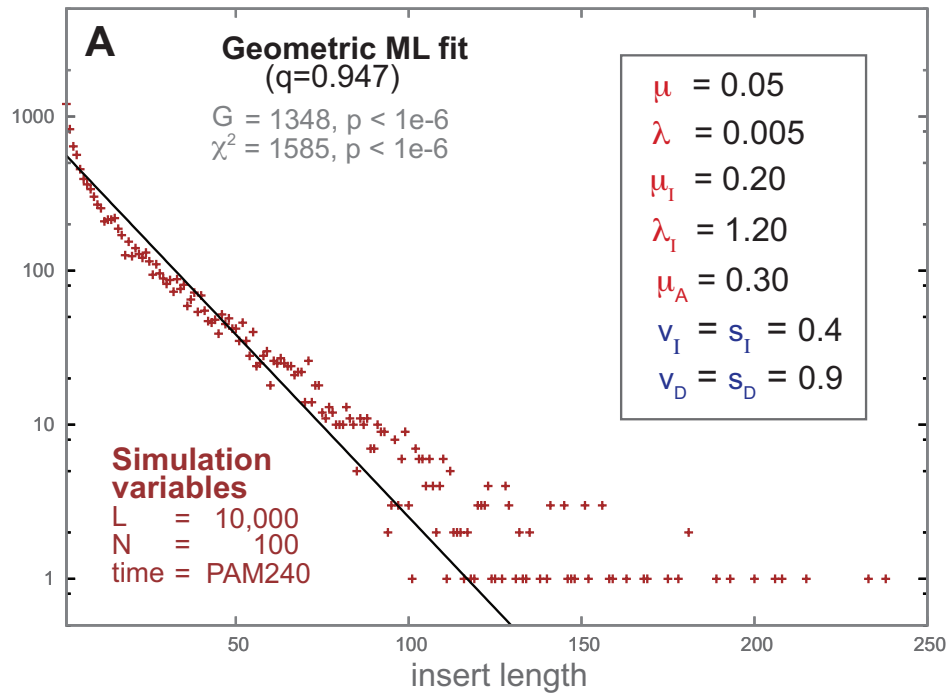
rate for adding to an insert "x" residue $\lambda_I (1 - v_I) v_I^{x-1}$

rate for removing from an insert "x" residues $\mu_I (1 - v_D) v_D^{x-1}$



$$P_t(\text{Descendant} | \text{Ancestral})$$

Not affine Models



Analytical closed-form solutions

AIF (fragment) Model

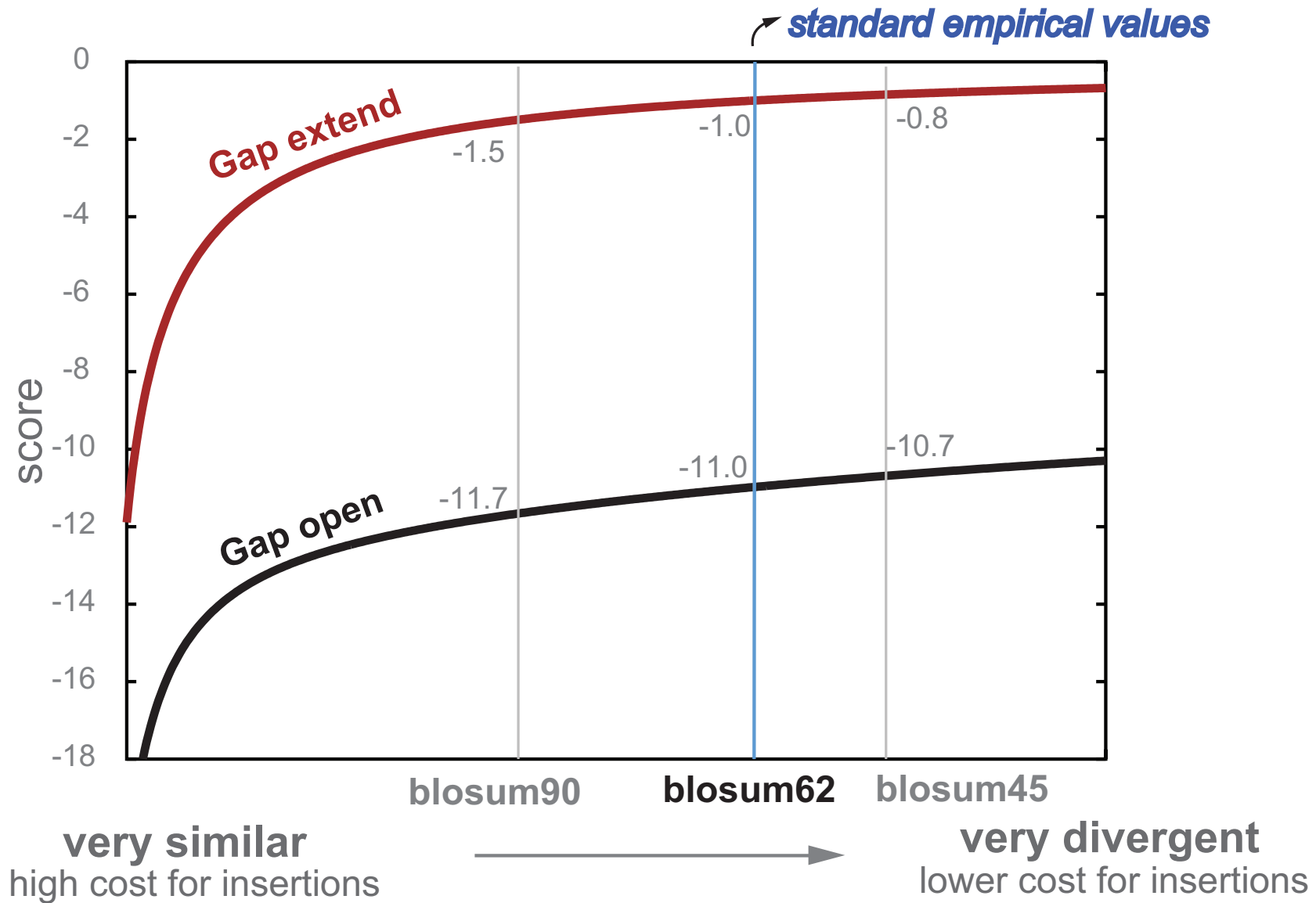
Gap opens: $\beta_t = \lambda_I \frac{1 - e^{(\lambda_I - \mu_I) t}}{\mu_I - \lambda_I e^{(\lambda_I - \mu_I) t}}$

Gap extends: $\eta_t = \frac{\lambda_I (1 - r) + \mu_I r - \lambda_I e^{(\lambda_I - \mu_I) t}}{\mu_I - \lambda_I e^{(\lambda_I - \mu_I) t}}$

**Ancestral
residue dies:** $\gamma_t = 1 - e^{-\mu_A t}$

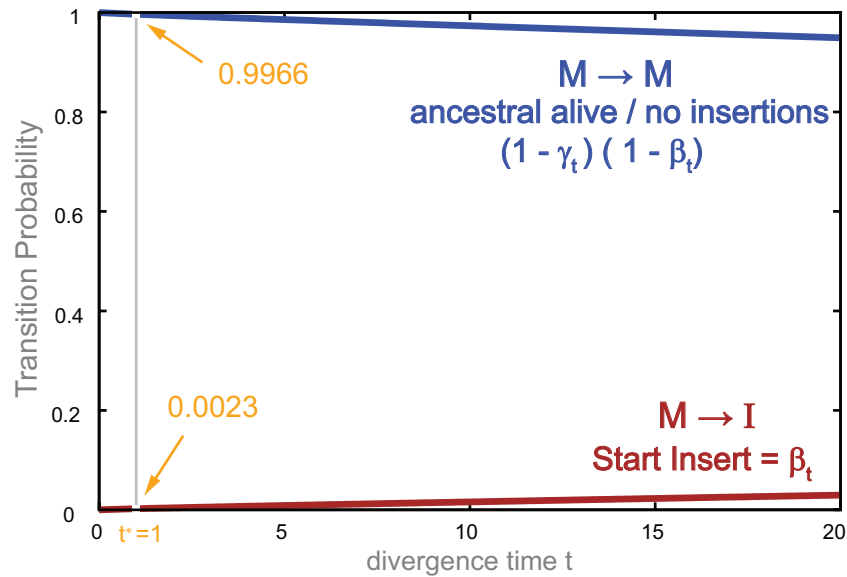
**More realistic microscopic models result
in non-affine macroscopic solutions**

Evolved BLAST

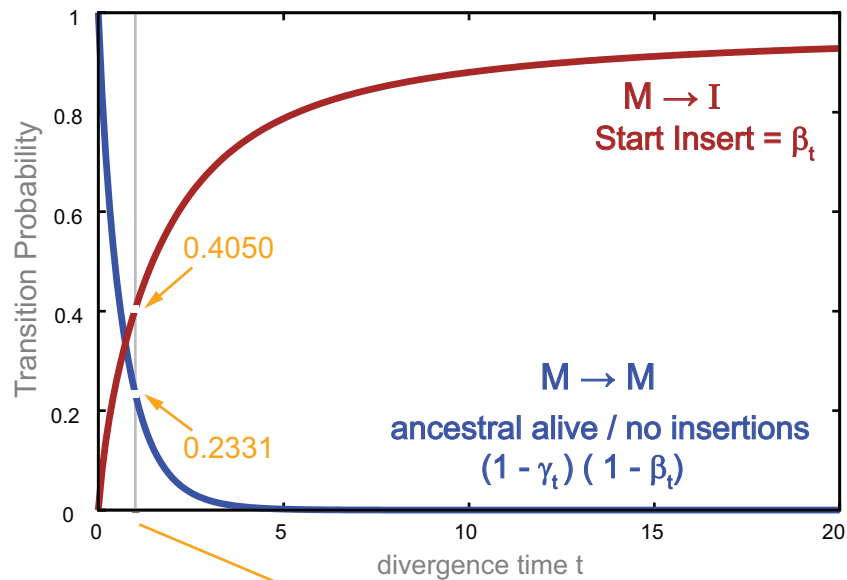


An evolved HMM

Position in a conserved region



Position at start of an insertion



time at which parameters were trained from data

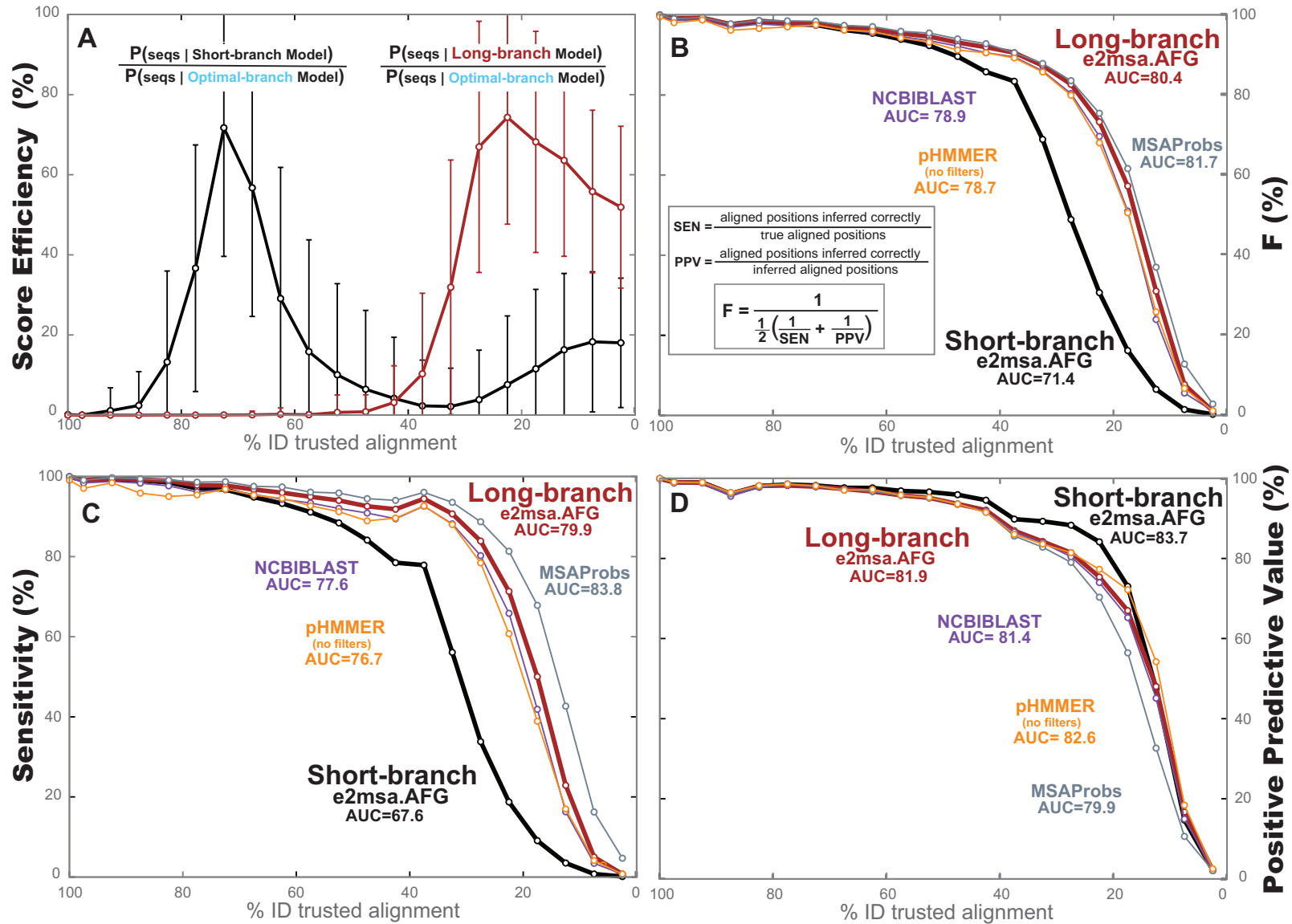
Affine Evolutionary Models

A Catalog

EVOLUTIONARY MODEL	total # free parameters	Microscopic model		Macroscopic model	
		rates	geometric parameters	# states minimal HMM	other properties
<i>single-residue models</i>					
AALI	6	$\lambda_I, \mu_I, \mu_A^{\{M,D,I\}}$	p	3	not reversible in general
LI	4	λ_I, μ_I, μ_A	p	1	not reversible in general
LR	2	$\lambda_I, \mu_A, (\mu_I = \lambda_I + \mu_A)$	$(p^{\text{LR}} = \lambda_I / \mu_A)$	1	reversible
TKF91	2	λ, μ	$(p^{\text{TKF}} = \lambda / \mu)$	2	reversible, ref. [?]
<i>fragment models</i>					
AFGX	9	$\lambda_I, \mu_I, \mu_A^{\{M,D,I\}}$	r_M, r_D, r_I, p	3	not reversible
AFG	7	λ_I, μ_I, μ_A	r_M, r_D, r_I, p	3	not reversible
AFGR	4	λ_I, μ_A	$r_M, r_{DI}, (p^{\text{LR}})$	3	reversible
AFR	3	λ_I, μ_A	$r, (p^{\text{LR}})$	3	reversible
TKF92	3	λ, μ	$r, (p^{\text{TKF92}})$	3	reversible, ref. [?]
FID	2	λ	$r, (p = 1)$	3	reversible, ref. [?]
<i>fragment affine model for profile HMMs</i>					
AIFX	7	$\lambda_I, \mu_I, \mu_A^{\{M,D,I\}}$	r_I, p	3	not reversible
AIF	5	λ_I, μ_I, μ_A	r_I, p	3	not reversible
<i>no-fragment affine model for profile HMMs (plan7 HMMER)</i>					
AGAX	9	$\lambda_{\{M,D\}}, \mu_{\{M,D\}}, \mu_A^{\{M,D,I\}}$	s_I, p	3	not reversible
AGA	7	$\lambda_{\{M,D\}}, \mu_{\{M,D\}}, \mu_A$	s_I, p	3	not reversible

A **fixed *long-branch***
parameterization is
sufficient to align
global homologies of all
degrees of
conservation.

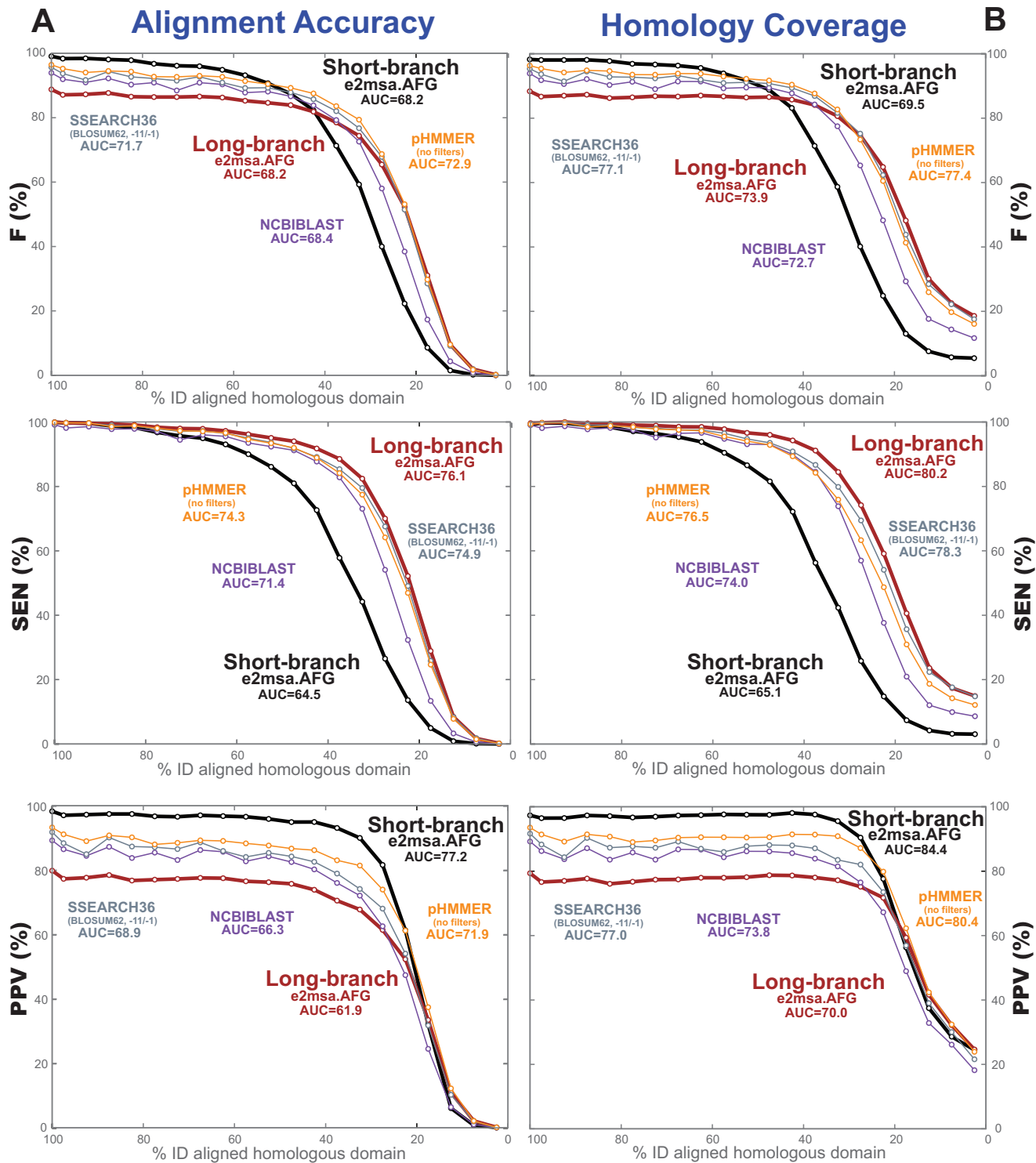
Global Homology Set



A **fixed *short-branch***
parameterization
reduces non-homologous
alignment overextension
for ***high-identity local***
homologies.

50 amino acid homologies

Local Homology Set

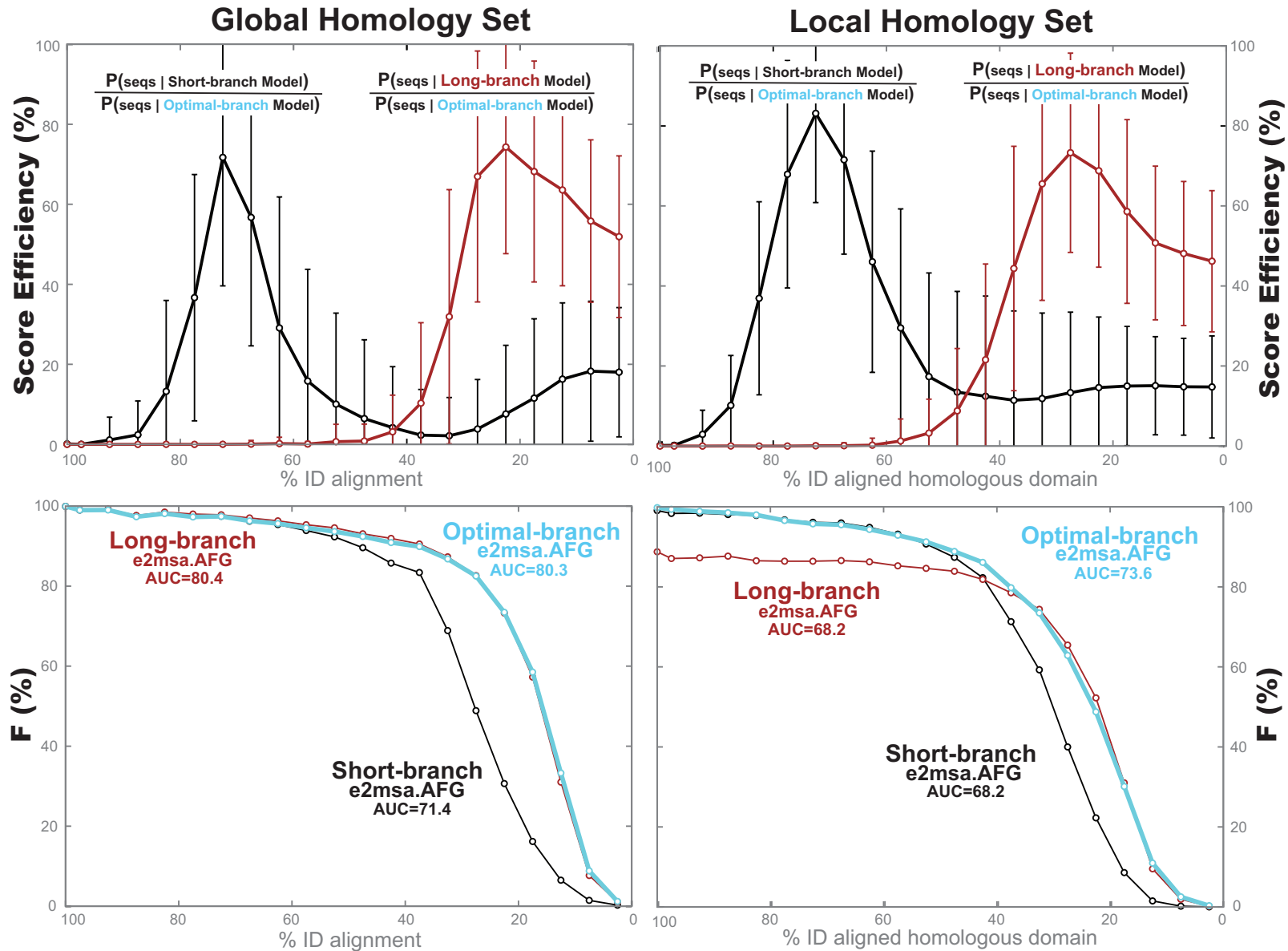


Optimal branch parameterization

A **variable *optimal-branch*** parameterization is best to align **local homologies** of any percentage identity.

e2msa - pairHMM aligner

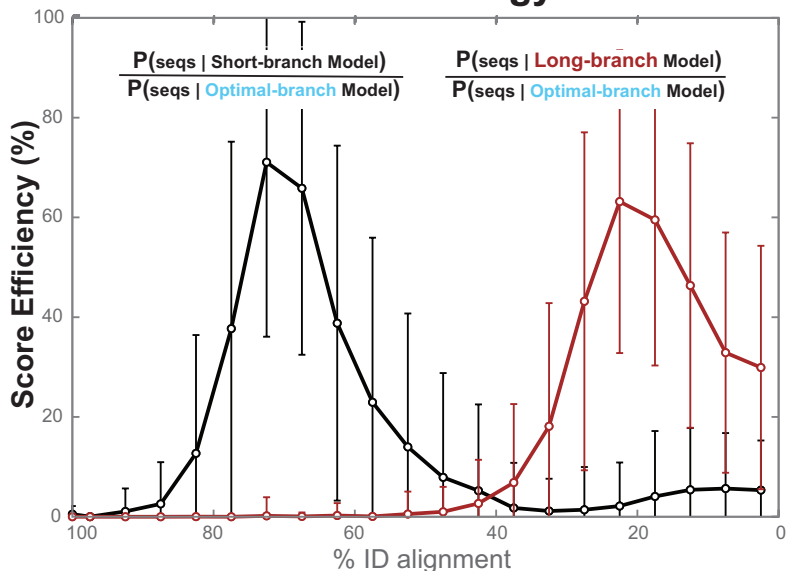
Alignment Accuracy- Evolutionary pair HMM (e2msa)



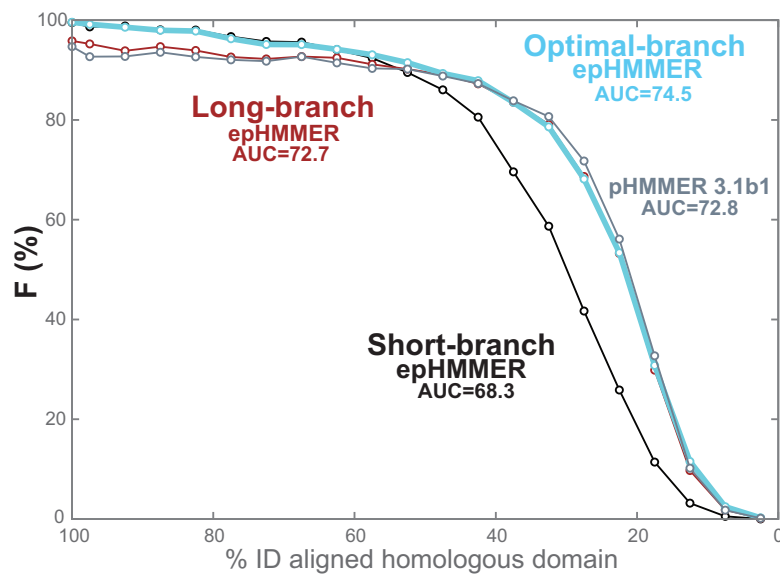
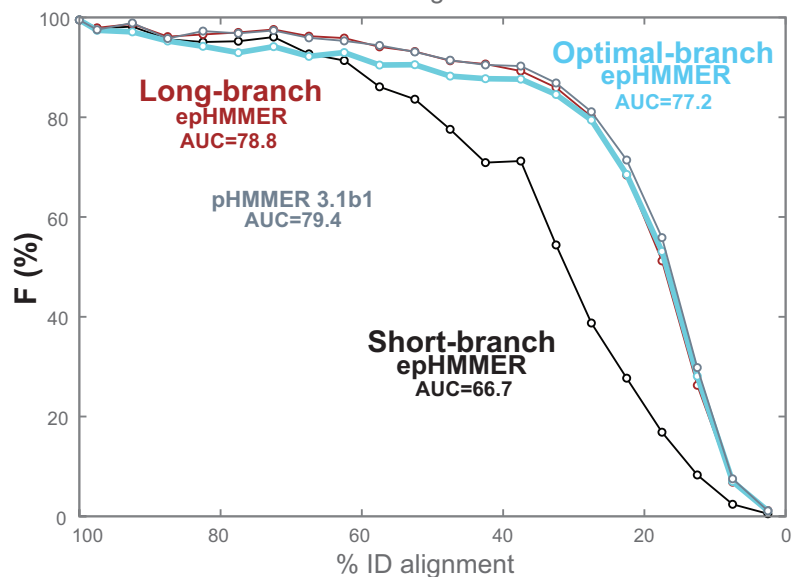
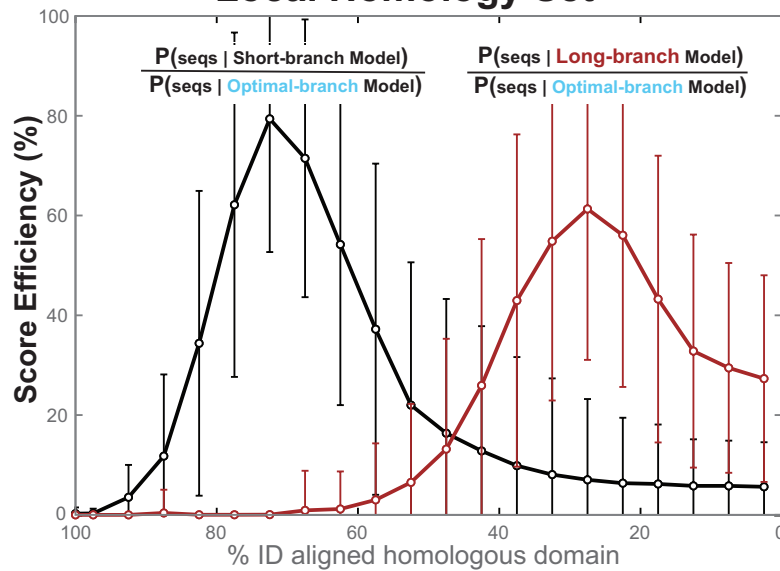
epHMMER

Alignment Accuracy - Evolutionary pHMMER (epHMMER)

Global Homology Set



Local Homology Set



Performance of different models

ALIGNMENT ACCURACY [AUC for F measure (%)]

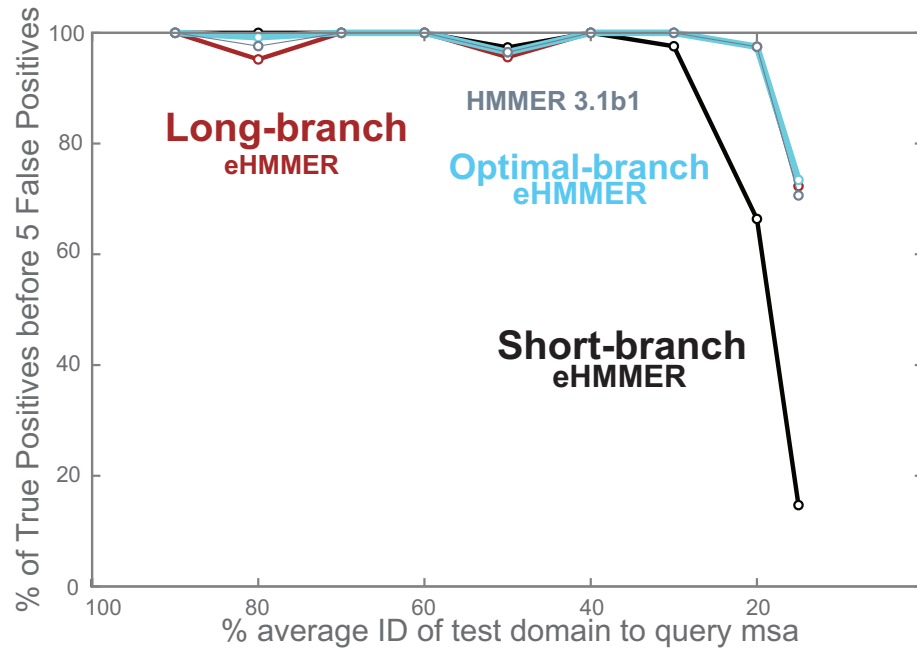
Method	Global Homology Set			Local Homology Set		
	PARAMETERIZATION			PARAMETERIZATION		
	SHORT	LONG	OPTIMAL	SHORT	LONG	OPTIMAL
e2msa.AFG	71.4	80.4	80.3	68.2	68.2	73.6
e2msa.AGA	71.3	80.4	80.2	68.1	67.3	73.6
e2msa.AIF	71.3	80.4	80.2	68.1	68.3	73.3
e2msa.TKF92	71.2	80.0	79.9	68.1	68.2	73.4
e2msa.LI	71.0	78.7	78.6	67.9	66.4	72.7
e2msa.TKF91	69.5	75.4	74.5	66.2	69.1	70.7
epHMMER (no filters)	66.7	78.8	77.2	68.3	72.7	74.5
pHMMER (no filters)		78.7			72.9	
SSEARCH (BLOSUM62, -11/-1)		80.0			71.7	
NCBIBLAST		78.9			68.4	
MSAProbs		81.7			36.1	
MUSCLE		80.8			33.5	

Evolutionary models with more parameters tend to perform better

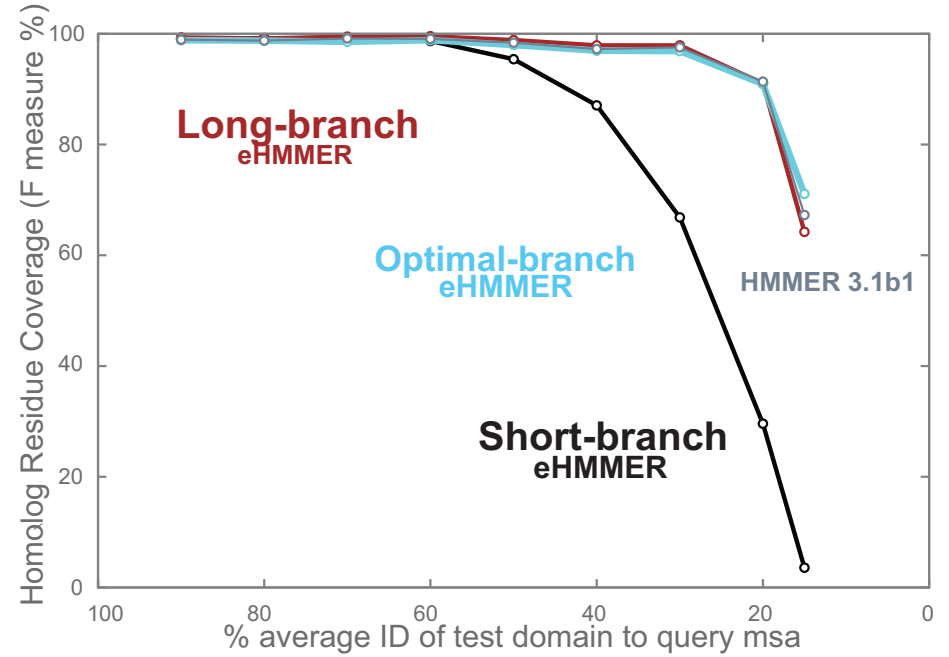
eHMMER

The *detection and coverage* of **embedded global homologies** is robust with just one **long-branch** parameterization

Embedded Global Homologies



Homology Detection

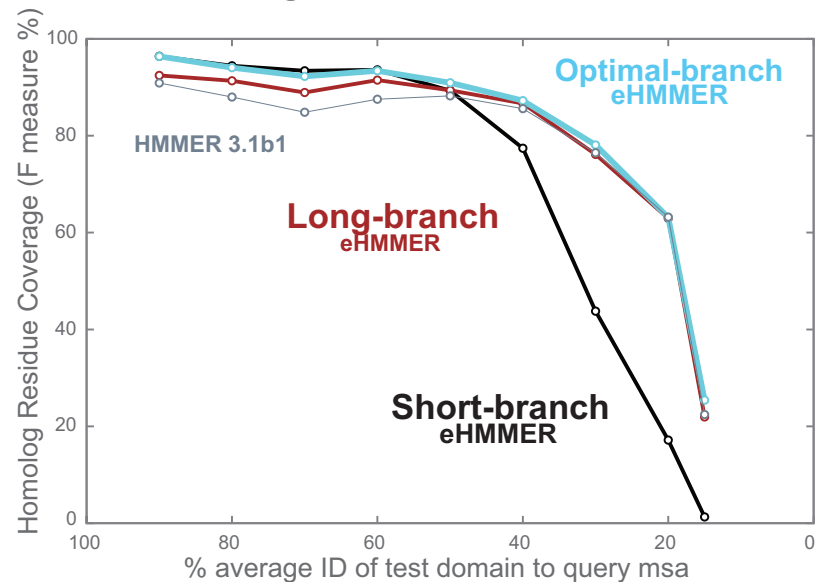
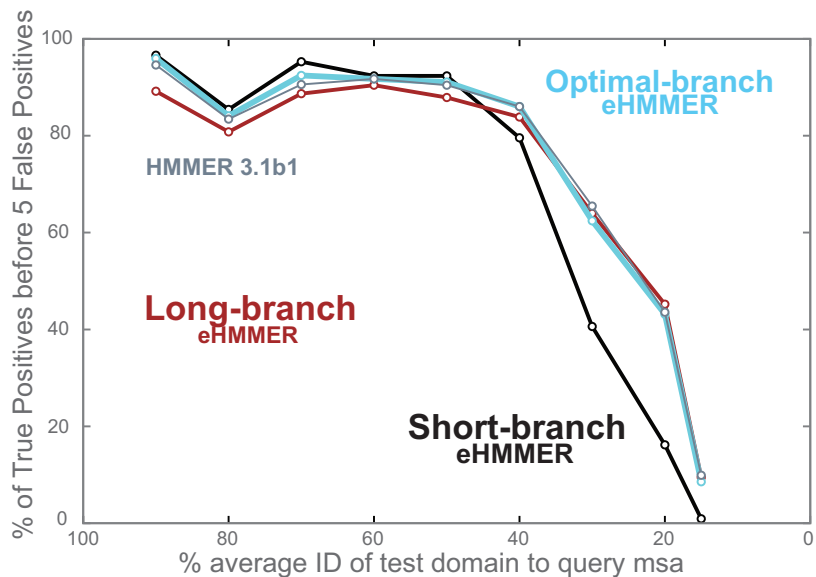


Homology Coverage

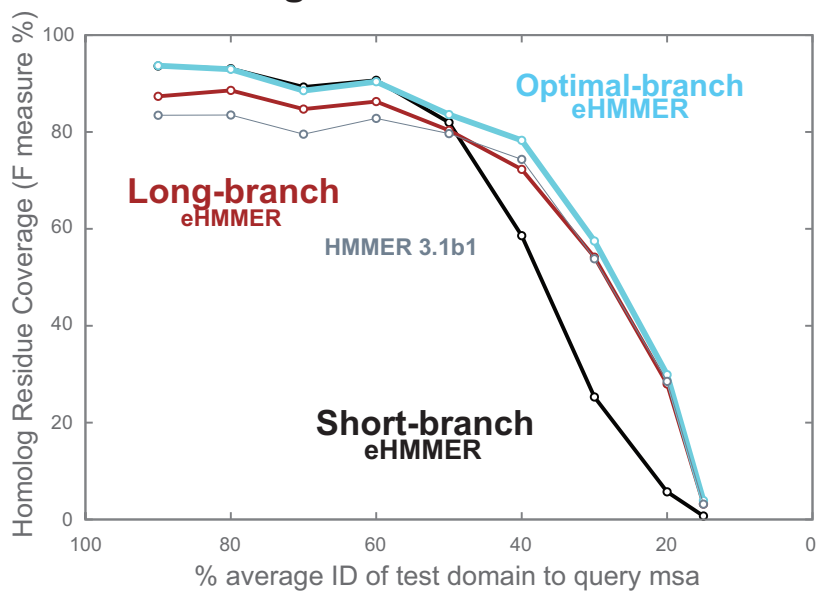
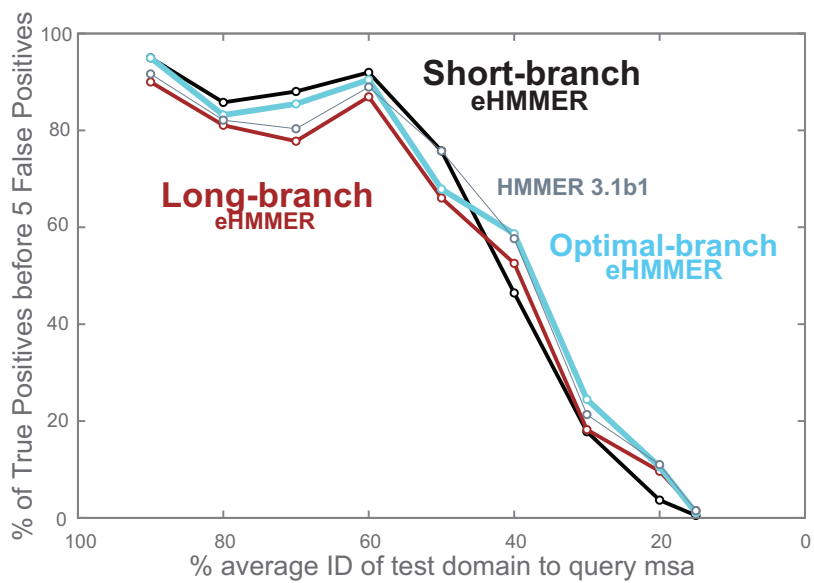
Short Local Homologies

The *detection and coverage* of **embedded short local homologies** improves with a variable **optimal-branch** parameterization

Embedded 50 aa Local Homologies



Embedded 30 aa Local Homologies



Homology Detection

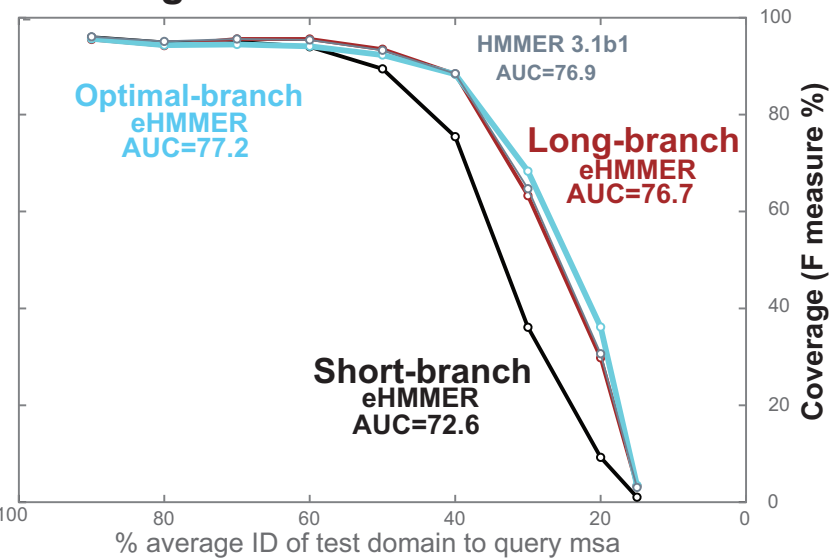
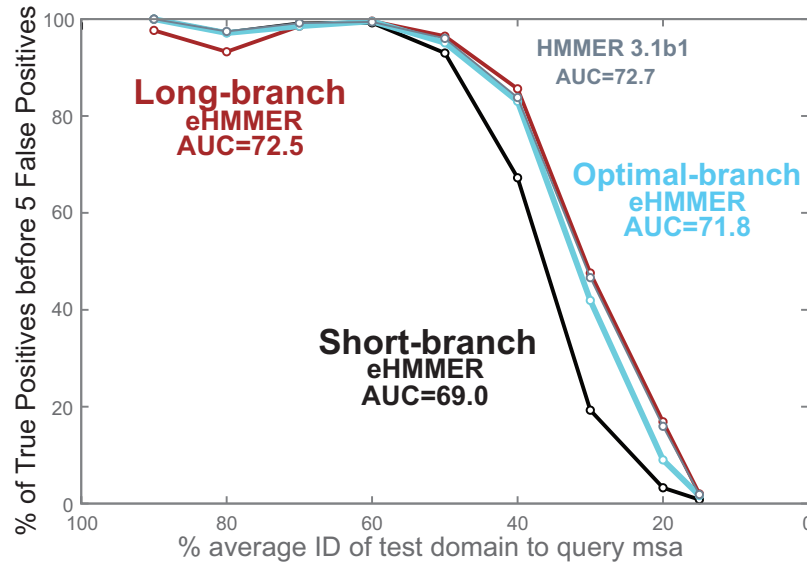
Homology Coverage

Fragments

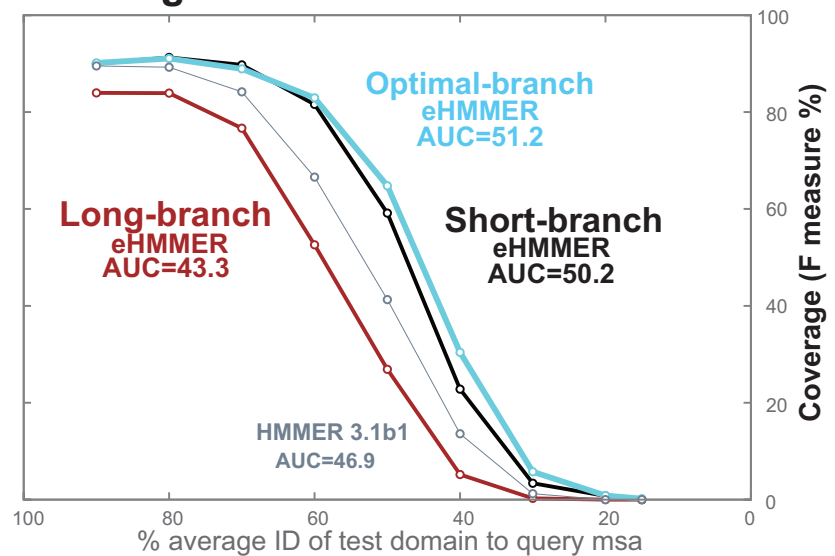
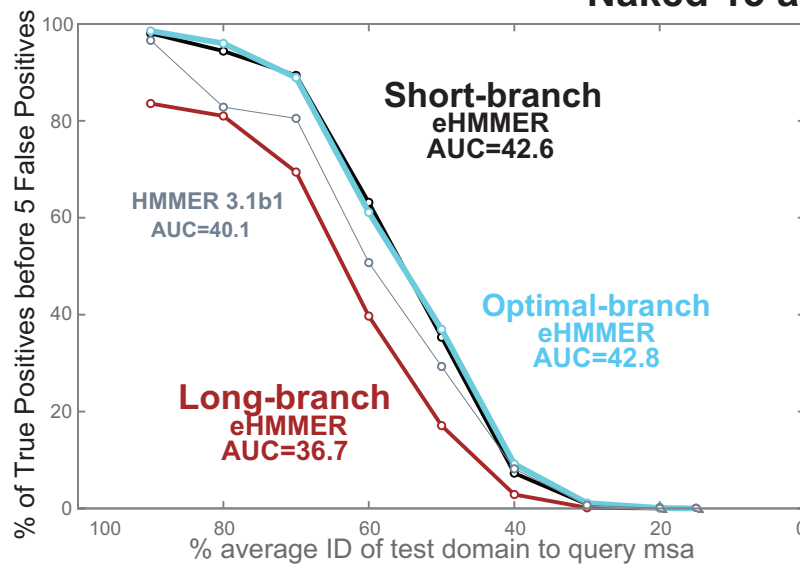
The *detection* of **very short** naked local homologies
improves with a
short-branch or optimal-branch
parameterization

Naked Fragments

Naked 30 aa Homologies



Naked 15 aa Homologies



Homology Detection

Homology Coverage

Explicit evolutionary models??

It is nice to wind up and down a model without additional information

- For Sensitivity → Use a **long-branch** parameterization (12% id).
Except for metagenomics < 30 aa,
then use a **short-branch** parameterization (45% id).
- For SEN/PPV → Use a **optimal-branch** for short embedded homologies.
For global embedded homologies still OK
using a **long-branch** parameterization.

Ancestral reconstruction

