

Boosting alignment accuracy through adaptive local realignment

Dan DeBlasio
John Kececioglu

Department of Computer Science
University of Arizona



Motivation

Multiple sequence alignment is a **fundamental problem** in bioinformatics.

- multiple sequence alignment is **NP-Complete**
- many **popular aligners** for multiple sequence alignment
- each aligner has many **parameters** whose values affect the alignment that is output

Motivation

Aligners often use *one* default **parameter choice** for *all* inputs.

- The **default** has good *average accuracy* across all benchmarks.
- The optimal default choice can be found by **inverse alignment** [Kececioglu and Kim 2007].
- The default may be a poor choice for **specific inputs**.

default

```
... yl-lhqflspssnqrtldqyggsvnrarlvlevvdavcnewsad-RIGIRVSPigtfq ...
... kP-LGVKLPPyf--dlvhfdimaeilnqfpltyvsnv-nsig---nglfidpeaesv ...
... yl-lnqfldphsnttrtdeyggnsienrarftlevvdalveaighe-KVGLRLSPygvfnd ...
... yl-plqflnpyynkrtldkyggslenrarfwletlekvkhavgsdcAIATRF---GVdt ...
... kvPLYVKLSPnv-tdivpiakaveaagadgltmintl-----mgvrfdlkrqp ...
```

alternate

```
... gsvenrarlvlevvdavcnewsad-RIGIRVSPigtfnvndngpnee--adalyl--- ...
... ydfeatekllke-----vftfftk-PLGVKLPPyf-----dlvhfdim ...
... gsienrarftlevvdalveaighe-KVGLRLSPygvfnsmsggaetgivaqyayvage ...
... gslenrarfwletlekvkhavgsdcAIATRFGV-----dtvygpgg ...
... tdpevaaalvka-----ckavskv-PLYVKLSPnvt-----divpiaka ...
```

Motivation

Proteins can have **different mutation rates** along their length

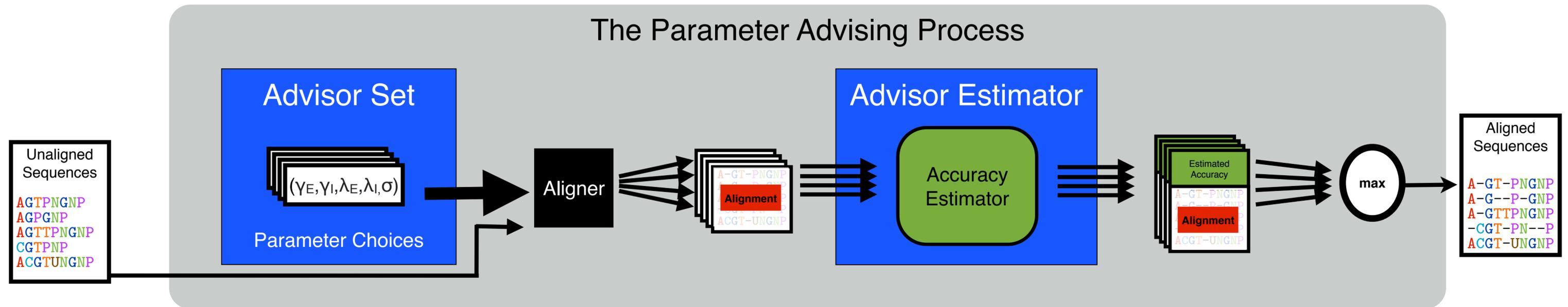
- Alignment parameters **model** mutation rates
- A **single choice** of parameters may not be best
- Using **different choices** across the protein can be superior

Can we find a parameter choice
that is best for *each region*
of a given input?

Parameter advising

Advising for unaligned input sequences

- aligns sequences using **each parameter choice** from a set,
- assigns an **estimated accuracy** to each alignment, and
- selects the alignment with the **highest estimated accuracy**.



Parameter advising

Advising for unaligned input sequences

- aligns sequences using **each parameter choice** from a set,
- assigns an **estimated accuracy** to each alignment, and
- selects the alignment with the **highest estimated accuracy**.



Accuracy estimation

Alignment accuracy is measured with respect to a reference alignment.

reference alignment	computed alignment
... a D E h s a D E h - s ...
... d S R - d d S R - - d ...
... a S H l t a S - H l t ...

- accuracy is the **fraction of substitutions** from the reference that are in the computed alignment,

Accuracy estimation

Alignment accuracy is measured with respect to a reference alignment.

reference alignment	computed alignment
... a D E h s a D E h - s ...
... d S R - d d S R - - d ...
... a S H l t a S - H l t ...

- accuracy is the **fraction of substitutions** from the reference that are in the computed alignment,

Accuracy estimation

Alignment accuracy is measured with respect to a reference alignment.

reference alignment	computed alignment
... a D E h s a D E h - s ...
... d S R - d d S R - - d ...
... a S H l t a S - H l t ...

- accuracy is the **fraction of substitutions** from the reference that are in the computed alignment,

Accuracy estimation

Alignment accuracy is measured with respect to a reference alignment.

reference alignment	computed alignment
... a D E h s a D E h - s ...
... d S R - d d S R - - d ...
... a S H l t a S - H l t ...
↑ ↑	

- accuracy is the **fraction of substitutions** from the reference that are in the computed alignment,
- measured on the **core columns** of the reference.

Accuracy estimation

Alignment accuracy is measured with respect to a reference alignment.

reference alignment	computed alignment	
... a D E h s a D E h - s ...	
... d S R - d d S R - - d ...	
... a S H l t a S - H l t ...	
↑ ↑		
		66% Accuracy

- accuracy is the **fraction of substitutions** from the reference that are in the computed alignment,
- measured on the **core columns** of the reference.

Accuracy estimation

Our estimator **Facet** (“**F**eature-based **AC**curacy **EsT**imator”)

- estimates accuracy by a **polynomial** on feature functions,
- efficiently learns the polynomial **coefficients** from examples,
- uses **novel features** that are efficient to evaluate.

Advisor sets

An **advisor set** should

- consist of **complementary** parameter settings,
- produce at least one **good alignment** for every input,
- be small, to reduce **running time**.

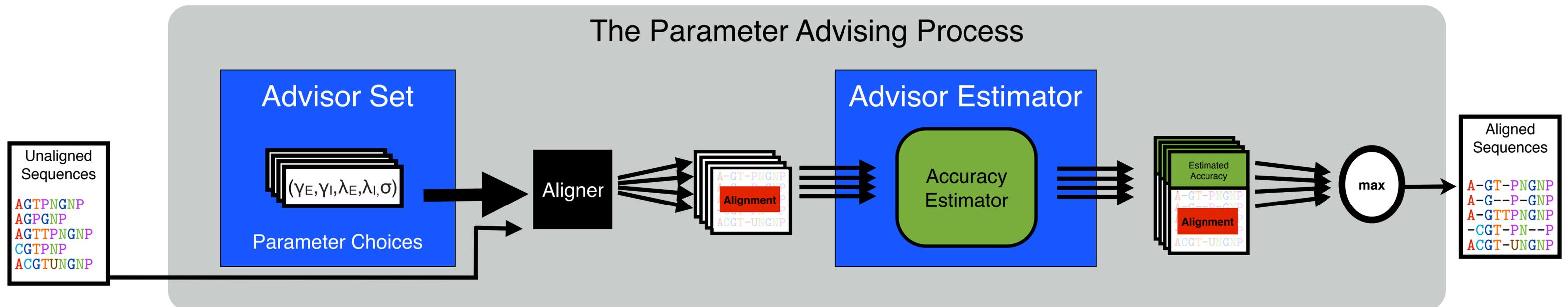
An **oracle set** is

- an **advisor set**, that is
- optimal for an **oracle advisor**,
- can be found by **integer linear programming**,
- can work well with an **actual advisor**.

Parameter advising

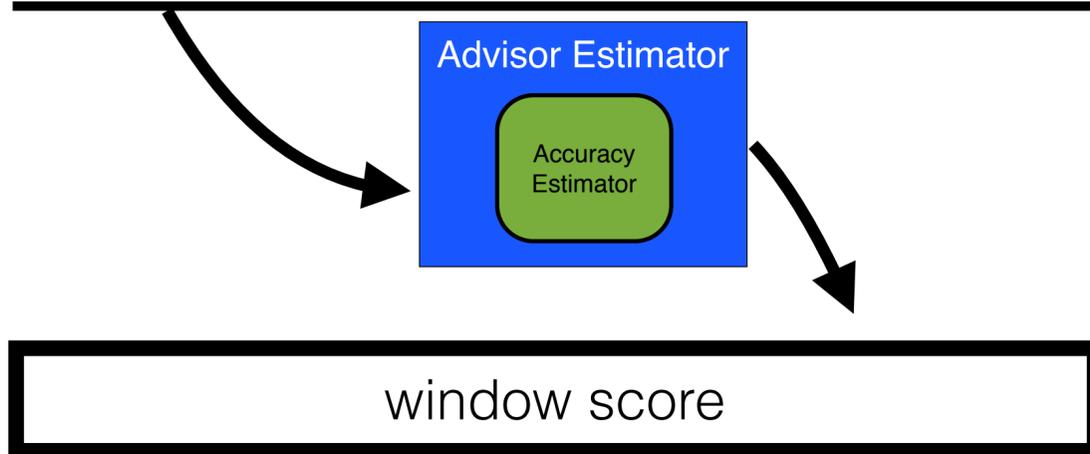
A **parameter advisor** has two components:

- an **accuracy estimator**, and
- a set of candidate **parameter choices**.



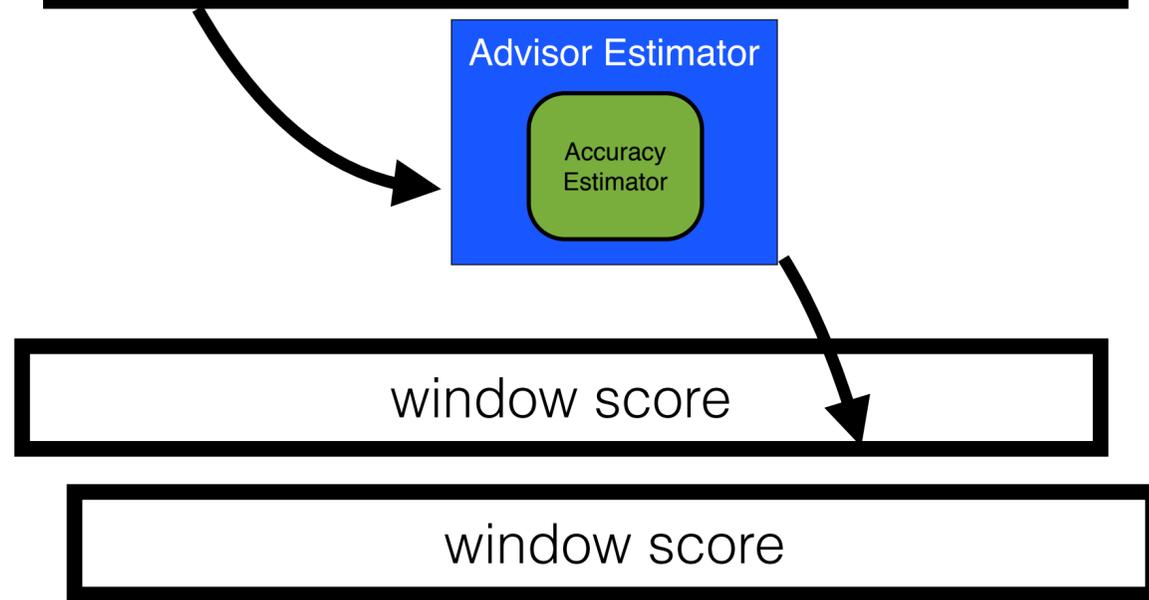
Adaptive local realignment

rkeyag**LYHEVAQA**HGVDVSQVrq**MKFGLEFFL**FDTLAVyenhfsnngvldqmsegrfafhkiindafttgychpnnd**PLVF**rwddsnaq**HKL**TLLVNQ**NDGEA**ARAEARVyleefvresysnt
kkaql**LYNEVATE**HGYDVTKId**MKFGNELL**FDTVWllehhftefgllldqmskgrfrfydlmkegfnegyaadne**PMIL**swiinthe**HCLSYITS**VDHDS**NR**AKD**ICRN**flghwy-dsyvna
rielln**HYQAAA**AKFNVDIANVr-----mtk**WNYGVFF**LYDVVA--**F**sehhdksyn**PLLF**kwddsqqk**HRLMLFVNV**NDNPT**QAKAEL**SIyledyl--sytqa
rlklls**FYNASAS**KYNKNIDLVR-----mnk**WNYGVFF**VYDVIN--**I**ddhylvkkds**PLVF**kwddinee**HQLMLHVNV**NEAETV**AKEELKL**yienyv--actqp



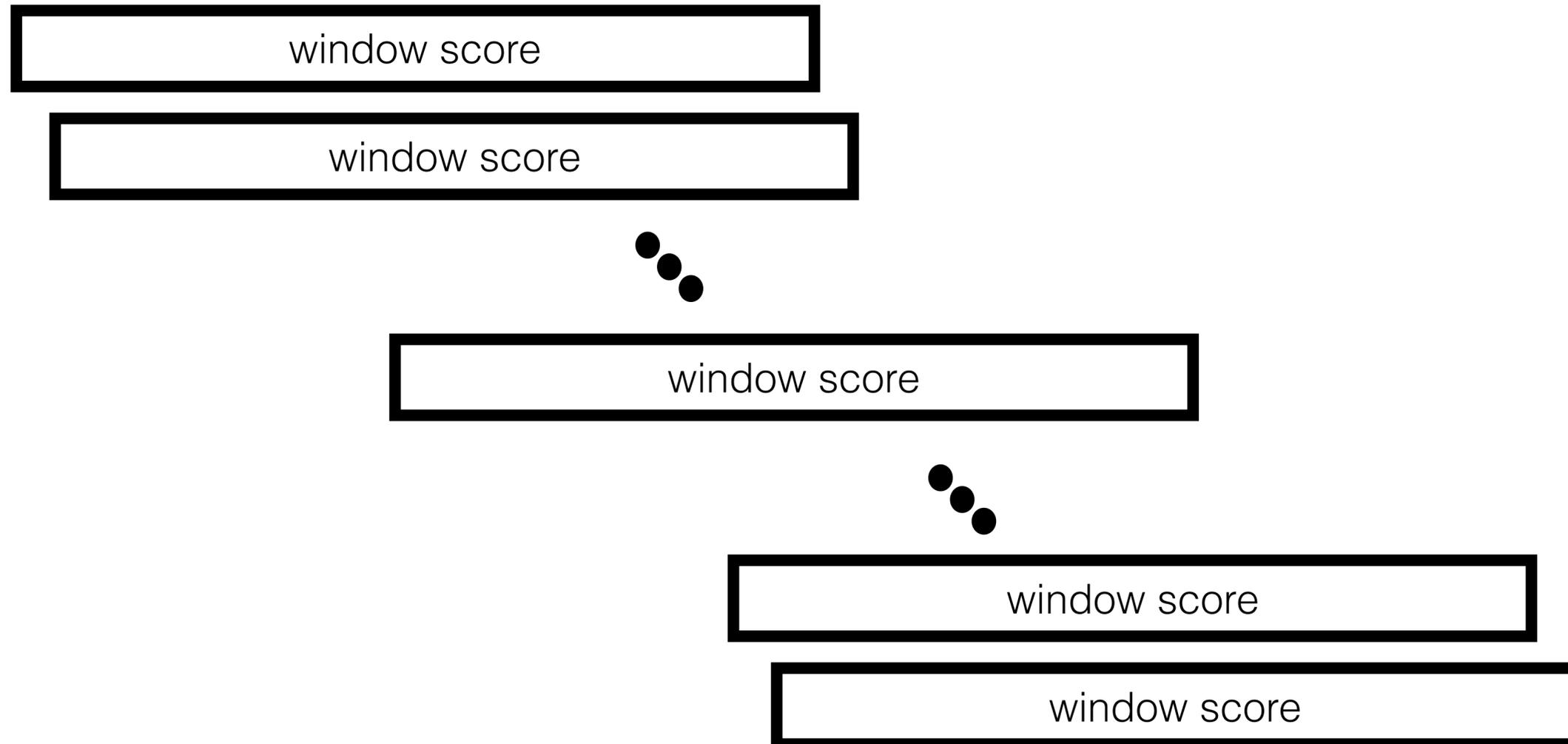
Adaptive local realignment

rkeyag**LYHEVAQA**HGVDVSQVrq**MKFGLEFL**FDTLAVyenhfsnngvldqmsegrfafhkiindafttgychpnnd**PLVF**rwddsnaq**HKL**TLLVNQ**NDGEA**ARAEARVyleefvresysnt
kkaql**LYNEVATE**HGYDVTKId**MKFGNELL**FDTVWllehhftefgllldqmskgrfrfydlmkegfnegyaadne**PMIL**swiinthe**HCLSYIT**SVDHDS**NRAKDI**CRNflghwy-dsyvna
rielln**HYQAAA**AKFNVDIANVr-----mtk**WNYGVFF**LYDVVA--**F**sehhdksyn**PLLF**kwddsqqk**HRLMLF**VNVNDNPT**QAKAEL**SIyledyl--sytqa
rlklls**FYNASAS**KYNKNIDLVR-----mnk**WNYGVFF**VYDVIN--**I**ddhylvkkds**PLVF**kwddinee**HQLMLH**VNVNEAET**VAKEEL**KLyienyv--actqp



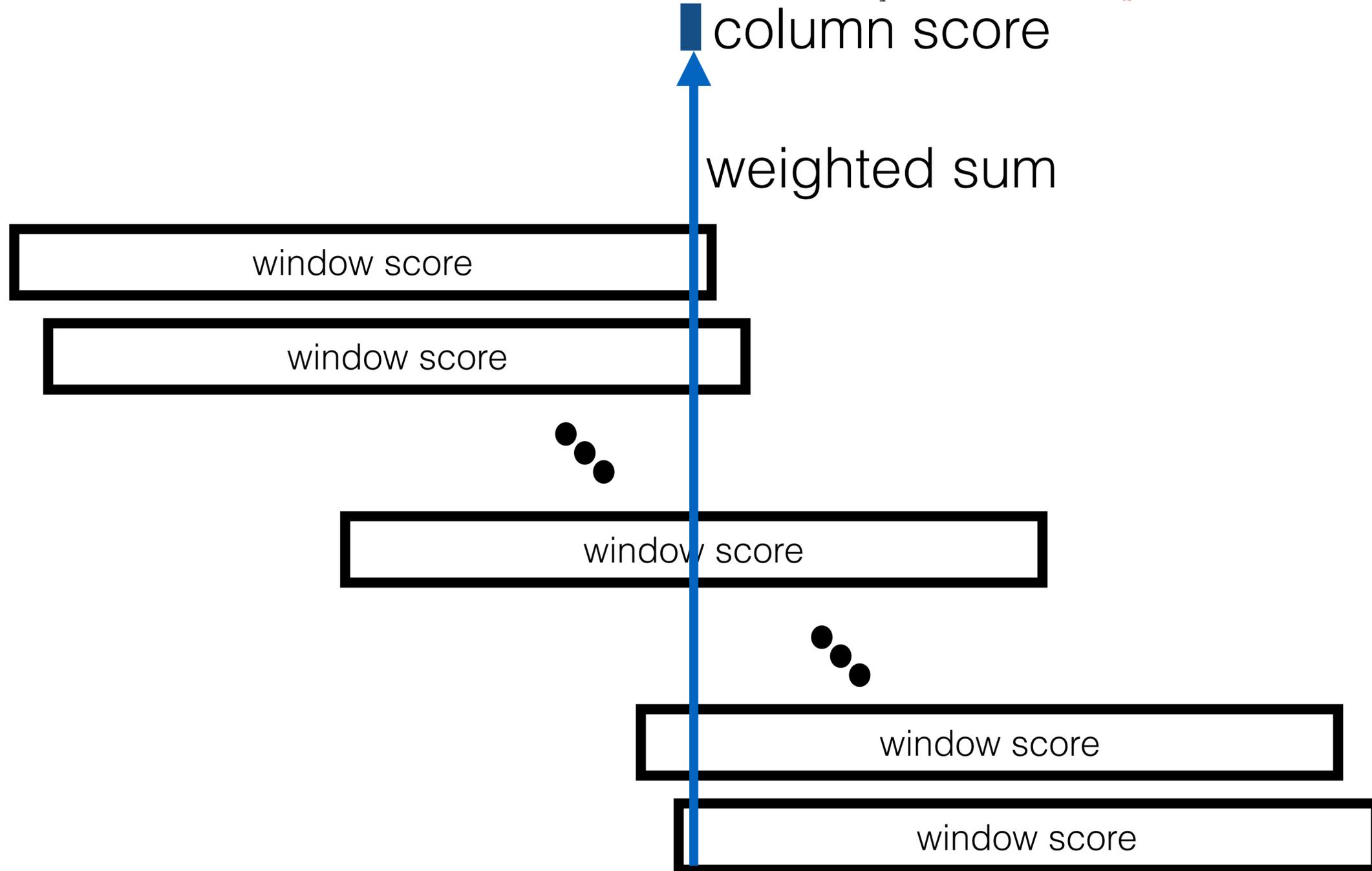
Adaptive local realignment

```
rkeyagLYHEVAQAHGVDVSQVrqMKFGLFFLFDTLAVyenhfsnngvldqmsegrfafhkiindafttgychpnndPLVFrwddsnaqHKLTLNQNVDGEAARAEARVyleefvresysnt  
kkaqlldLYNEVATEHGVDVTKId-MKFGNELLFDVWVWlehhftefgllldqmskgrfrfydlmkegfnegyaadnePMILswiintheHCLSYITSVDHDSNRAKDICRNflghwy-dsyvna  
riellnHYQAAAKFNVDIANVr-----mtkWNYGVFFLYDVVA--FsehhidksynPLLFkwddsqqkHRLMLFVNVNDNPTQAKAELSIyledyl--sytqa  
rklslsFYNASASKYNKNIDLVR-----mnkWNYGVFFVYDVIN--IddhylvkkdsPLVFKwddineeHQLMLHVNVNEAETVAKEELKLyienyv--actqp
```



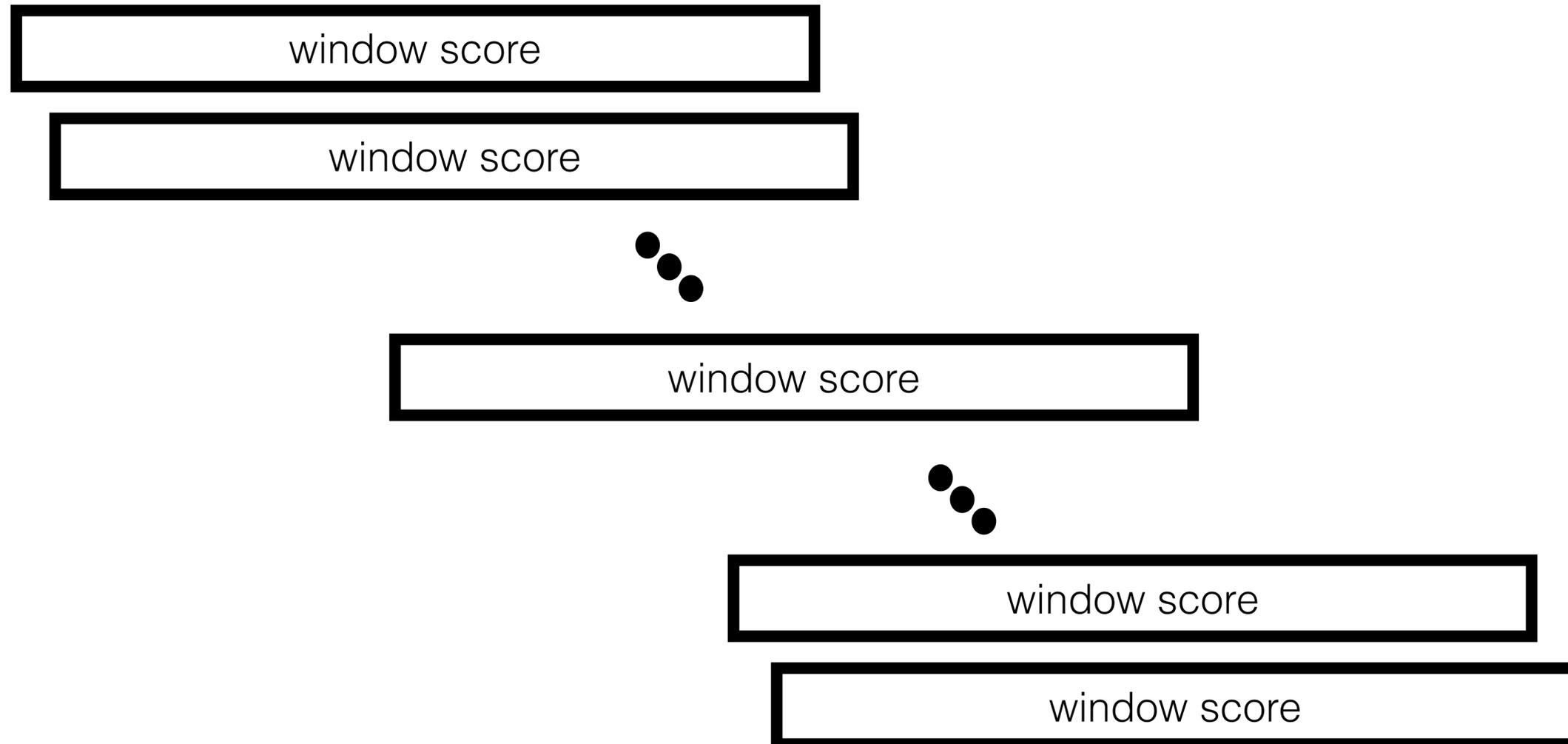
Adaptive local realignment

```
rkeyagLYHEVAQAHGVDVSQVrqMKFGLFFLFDTLAVyenhfsnngvldqmsegrfafhkiindafttgychpnndPLVFrwddsnaqHKLTLNQNVDGEAARAEARVyleefvresysnt  
kkaqldLYNEVATEHGVDVTKId-MKFGNELLFDVWVWlehhftefgllldqmskgrfrfydlmkegfnegyaadnePMILswiintheHCLSYITSVDHDSNRAKDICRNflghwy-dsyvna  
riellnHYQAAAKFNVDIANVr-----mtkWNYGVFFLYDVVA--FsehhdksynPLLFkwddsqqkHRLMLFVNVNDNPTQAKAELSIyledyl--sytqa  
rklslsFYNASASKYNKNIDLVR-----mnkWNYGVFFVYDVIN--IddhylvkkdsPLVFkwddineehQLMLHVNVEAETVAKEELKLyienyv--actqp
```



Adaptive local realignment

```
rkeyagLYHEVAQAHGVDVSQVrqMKFGLFFLFDTLAVyenhfsnngvldqmsegrfafhkiindafttgychpnndPLVFrwddsnaqHKLTLLVNQNVGDGEAARAEARVyleefvresysnt  
kkaqldLYNEVATEHGVDVTKId-MKFGNELLFDTVWLlehhftefgllldqmskgrfrfydlmkegfnegyaadnePMILswiintheHCLSYITSVDHDSNRAKDICRNflghwy-dsyvna  
riellnHYQAAAKFNVDIANVr-----mtkWNYGVFFLYDVVA--FsehhidksynPLLFkwddsqqkHRLMLFVNVNDNPTQAKAELSIyledyl--sytqa  
rlkllsFYNASASKYNKNIDLVR-----mnkWNYGVFFVYDVIN--IddhylvkkdsPLVfkddineeHQLMLHVNVNEAETVAKEELKLyienyv--actqp
```



Adaptive local realignment



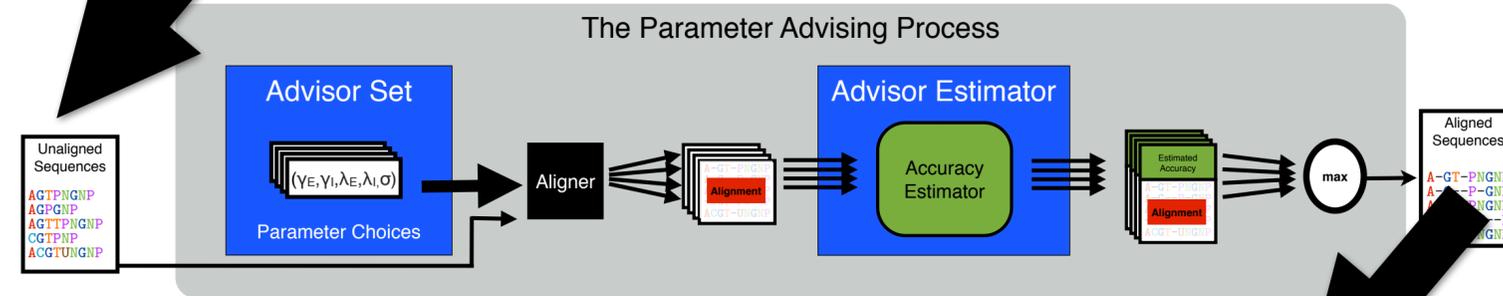
Adaptive local realignment

```

rkeyagLYHEVAQAHGVDVSQVr--qMKFGLFFLEFDTLAVyenhfsnngvldqmsegrfafhkiindafttgychpnndELVFrwddsnaqHKLTLNQNVDGEAARAEARVyleefvresysnt
kkaqldLYNEVATEHGYDVTKIc---MKFGNFFLEFDTVWLlehhftefgllldqmskgrfrfydlmkegfnegyaadneEMILswiintheHCLSYITSDHDSNRAKDCRNflghwy-dsyvna
riellnHYQAAAANKFNVDIANVr-----mtkWNYGVFVFLYDVVA--FsehhidksynELLFkwddsqqkHRLMLFVNVDNPTQAKAELSIyledyl--sytqa
rlkllsFYNASASKYNKNIDLVR-----mnkWNYGVFVYDVVIN--IddhylvkkdsELVFKwddineeHQLMLHVNNEAETVAKEELKLyienyv--actqp
    
```

```

qMKFGLFFLEFDTLAVyenhfsnngvldqmsegrfafhkiindafttgychpnnd
--MKFGNFFLEFDTVWLlehhftefgllldqmskgrfrfydlmkegfnegyaadne
-----mtkWNYGVFVFLYDVVA--Fsehhidksyn
-----mnkWNYGVFVYDVVIN--Iddhylvkkds
    
```



```

--qMKFGLFFLEFDTLAVyenhfsnngvldqmsegrfafhkiindafttgychpnnd
---MKFGNFFLEFDTVWLlehhftefgllldqmskgrfrfydlmkegfnegyaadne
mtkWNYGVFVFLYDVVAFsehhidksyn-----
mnkWNYGVFVYDVVINIddhylvkkds-----
    
```

```

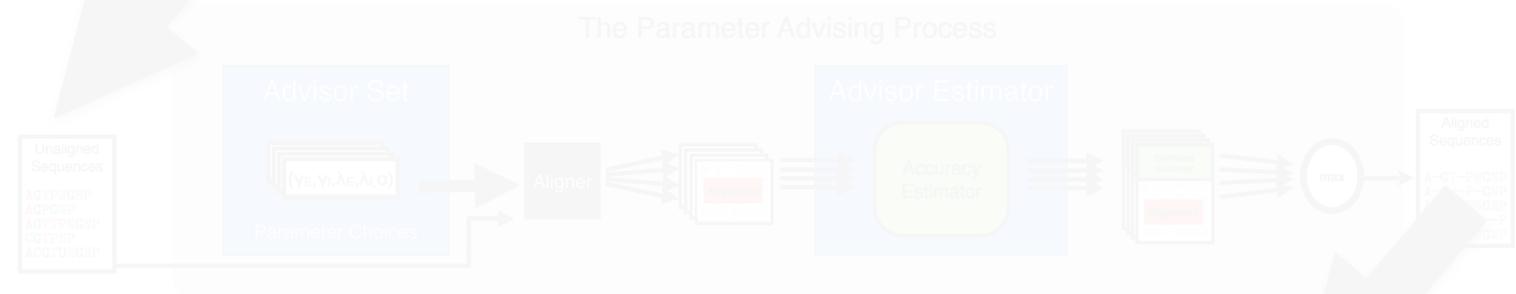
rkeyagLYHEVAQAHGVDVSQVr--qMKFGLFFLEFDTLAVyenhfsnngvldqmsegrfafhkiindafttgychpnndELVFrwddsnaqHKLTLNQNVDGEAARAEARVyleefvresysnt
kkaqldLYNEVATEHGYDVTKIc---MKFGNFFLEFDTVWLlehhftefgllldqmskgrfrfydlmkegfnegyaadneEMILswiintheHCLSYITSDHDSNRAKDCRNflghwy-dsyvna
riellnHYQAAAANKFNVDIANVr-----mtkWNYGVFVFLYDVVAFsehhidksyn-----ELLFkwddsqqkHRLMLFVNVDNPTQAKAELSIyledyl--sytqa
rlkllsFYNASASKYNKNIDLVR-----mnkWNYGVFVYDVVINIddhylvkkds-----ELVFKwddineeHQLMLHVNNEAETVAKEELKLyienyv--actqp
    
```

Adaptive local realignment

```
rkeyagLYHEVAQAHGVDVSQVrMKFGLFFLFDTLAVyenhfsnngvldqndFLVFrwddsnaqHKLTLNVNQVNDGEAARAEARVyleefvresysnt  
kkaqldLYNEVATEHGYDVTKIdMKFGNFLFLFDTVWllehhftefgllldqneMILswiintheHCLSYITSVDHDSNRKADICRNflghwy-dsyvna  
riellnHYQAAAANKFNVDIANVr-----mtk-----ynELLFkwddsqqkHRLMLFVNVNDNPTQAKAELSIyledyl--sytqa  
rklslsFYNASASKYNKNIDLVR-----mnk-----dsFLVfkwdineehQQLMLHVVNVNEAETVAKEELKLyienyv--actqp
```

Accuracy
33%

```
qMKFGLFFLFDTLAVyenhfsnngvldqmsegrfafhkiindafttgychpnnd  
-MKFGNFLFLFDTVWllehhftefgllldqmskgrfrfydlmkegfnegyiaadne  
-----mtkWNYGVFFLYDVVA--Fsehhidksyn  
-----mnkWNYGVFFVYDVINI--Iddhylvkkds
```



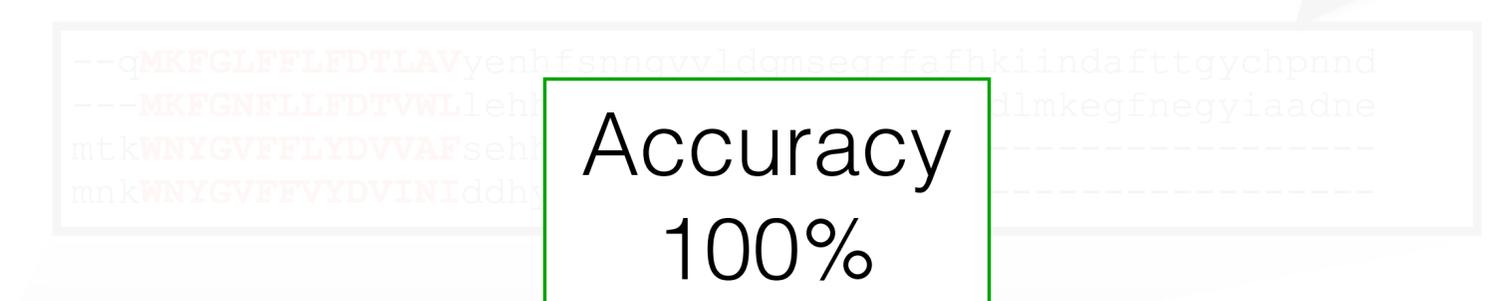
```
--qMKFGLFFLFDTLAVyenhfsnngvldqmsegrfafhkiindafttgychpnnd  
---MKFGNFLFLFDTVWllehhftefgllldqmskgrfrfydlmkegfnegyiaadne  
mtkWNYGVFFLYDVVAFsehhidksyn-----  
mnkWNYGVFFVYDVINIddhylvkkds-----
```

```
rkeyagLYHEVAQAHGVDVSQVr--qMKFGLFFLFDTLAVyenhfsnngvldqndFLVFrwddsnaqHKLTLNVNQVNDGEAARAEARVyleefvresysnt  
kkaqldLYNEVATEHGYDVTKId---MKFGNFLFLFDTVWllehhftefgllldqneMILswiintheHCLSYITSVDHDSNRKADICRNflghwy-dsyvna  
riellnHYQAAAANKFNVDIANVrmtkWNYGVFFLYDVVAFsehhidksyn---ELLFkwddsqqkHRLMLFVNVNDNPTQAKAELSIyledyl--sytqa  
rklslsFYNASASKYNKNIDLVRmnkWNYGVFFVYDVINIddhylvkkds---FLVfkwdineehQQLMLHVVNVNEAETVAKEELKLyienyv--actqp
```

Accuracy
100%

Adaptive local realignment

```
rkeyagLYHEVAQAHGVDVSQVrMKFGLFFLFDTLAVyenhfsnngvldqmsegrfafhkiindafttgychpnndELVFrwddsnaqHKLTLNVNQNVDGEAARAEARVyleefvresysnt  
kkaqldLYNEVATEHGYDVTKI---MKFGNLLFDTVWllehhftefgllldqmskgrfrfydlmkegfnegyaadneEMILswiintheHCLSYITSVDHDSNRAKDICRNflghwy-dsyvna  
riellnHYQAAAKFNVDIANVr-----mtkWNYGVFFLYDVVA--FsehhidksynELLFkwddsqqkHRLMLFVNVDNPTQAKAELSIyledyl--sytqa  
rlkllsFYNASASKYNKNIDLv-----mnkWNYGVFFVYDVIN--IddhylvkkdsELVFWkddineeHQLMLHVNVEAETVAKEELKLyienyv--actqp
```



```
rkeyagLYHEVAQAHGVDVSQVr--qMKFGLFFLFDTLAVyenhfsnngvldqmsegrfafhkiindafttgychpnndELVFrwddsnaqHKLTLNVNQNVDGEAARAEARVyleefvresysnt  
kkaqldLYNEVATEHGYDVTKI---MKFGNLLFDTVWllehhftefgllldqmskgrfrfydlmkegfnegyaadneEMILswiintheHCLSYITSVDHDSNRAKDICRNflghwy-dsyvna  
riellnHYQAAAKFNVDIANVrmtkWNYGVFFLYDVVAFsehhidksyn-----ELLFkwddsqqkHRLMLFVNVDNPTQAKAELSIyledyl--sytqa  
rlkllsFYNASASKYNKNIDLvmnkWNYGVFFVYDVINIddhylvkkds-----ELVFWkddineeHQLMLHVNVEAETVAKEELKLyienyv--actqp
```

Adaptive local realignment

Adaptive local realignment takes as **input**

- an **initial alignment**,
- an **advisor set**,
- window **sizes**, and
- column score **thresholds**.

And **outputs**

- a **new alignment**
- constructed by **realigning** regions of low estimated accuracy.

Experimental results

We **evaluate** the accuracy of adaptive local realignment

- with the Opa1 **aligner** and Facet **estimator**,
- on over 800 **benchmarks** from BENCH and PALI,
- using 12-fold **cross-validation**.

Experimental results

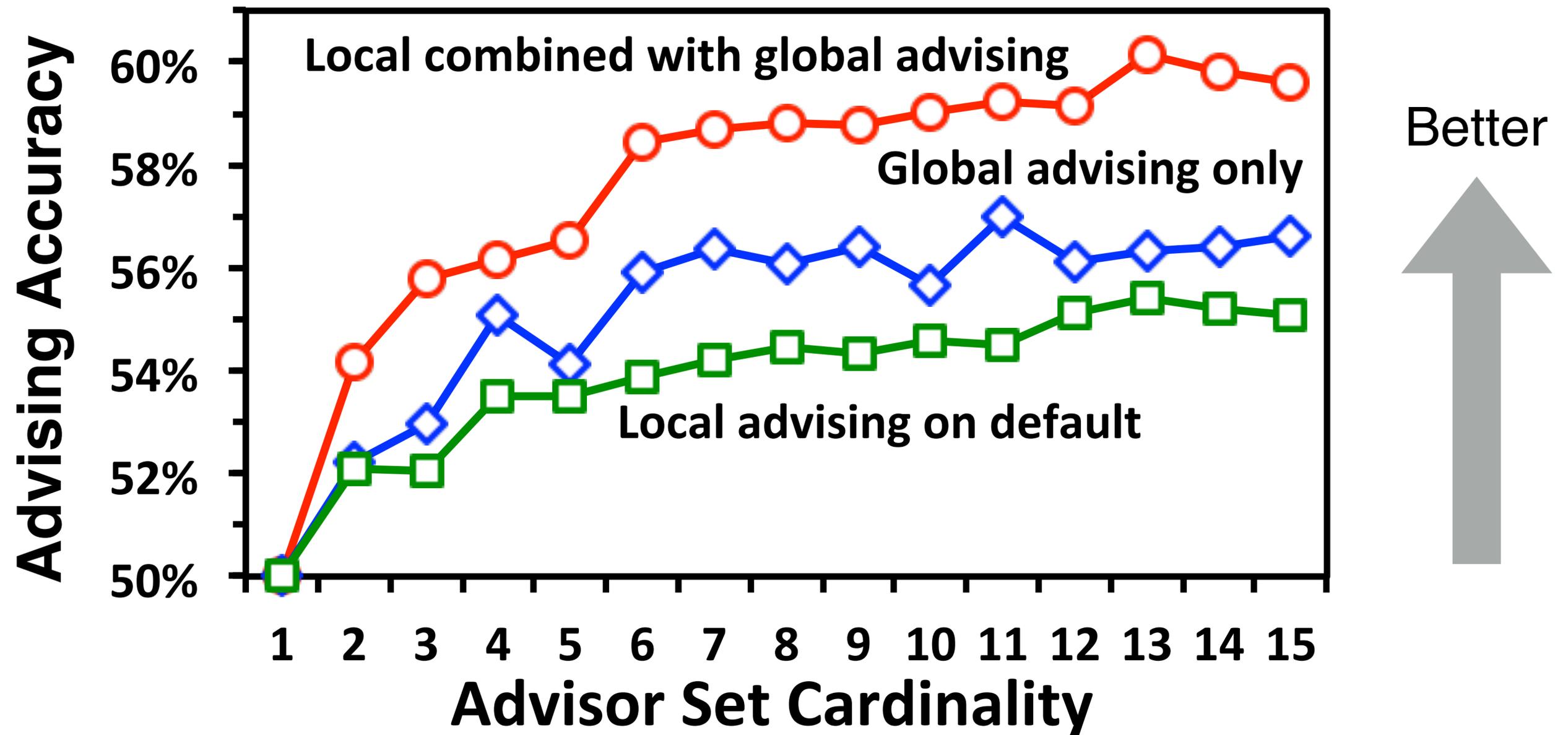
We correct for the **bias** in over-representation of easy-to-align benchmarks.

- The **difficulty** of a benchmark is its accuracy under the default parameter setting.
- Split the range of difficulties $[0,1]$ into **10 bins**.
- Report advisor accuracy uniformly **averaged** across bins.

The typical **average accuracy** is close to 50%.

Experimental results

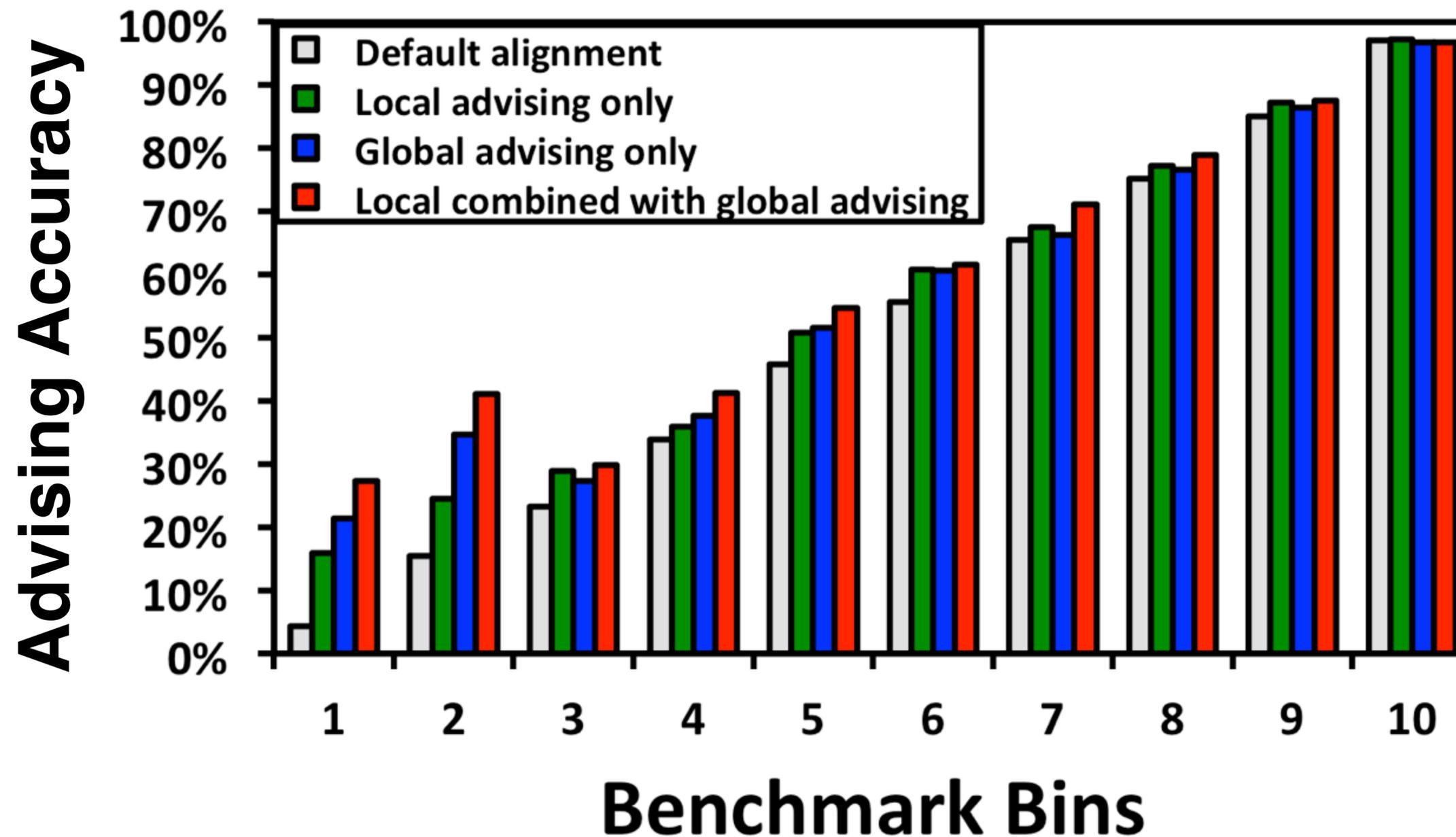
Average accuracy of advisors versus set cardinality



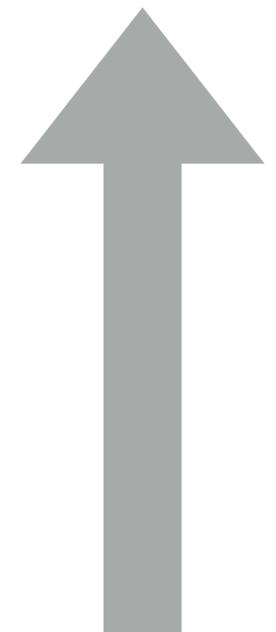
Boosts accuracy by 9% over default parameter choice

Experimental results

Average accuracy of advisors within difficulty bins



Better



oracle advising set
 $k = 10$

Boosts accuracy on the hardest bins by almost 26%

Conclusions

Local adaptive realignment:

- Uses **diverse** parameter choices for heterogeneous proteins.
- Facet estimator identifies misaligned **regions**.
- Misaligned regions are polished with **parameter advising**.
- With global advising, **boosts accuracy** by **up to 26%**.

Further research

Promising directions include:

- Finding **advisor sets** tuned for **local realignment**.
- Improving the **estimator** by training on **local examples**.
- Using an **ensemble of aligners** to **realign regions**.

Software distribution

Available for download:

- **Facet** accuracy estimator
- **Opal** aligner with **global** and **local** parameter advising
- **Parameter sets** for advising

facet.cs.arizona.edu

Acknowledgments

People

William Pearson

Travis Wheeler

Poster

N22

Monday, July 11

Funding

- University of Arizona
NSF IGERT in Genomics
Grant DGE-0654435
- NSF Grant IIS-1217886

