

Boosting alignment accuracy through adaptive local realignment

Dan DeBlasio

Computational Biology Department
Carnegie Mellon University*

John Kececioglu

Department of Computer Science
University of Arizona

Motivation

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

Motivation

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

- multiple sequence alignment is **NP-Complete**

Motivation

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

- multiple sequence alignment is **NP-Complete**
- many **popular aligners** for multiple sequence alignment

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.

Motivation

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

- The **default** has good *average accuracy* across all benchmarks.

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**.

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**.
- The default may be a poor choice for **specific inputs**.

	... yl-1hqflspssnqtdqyggsvenrarlvlevvdavcnewsad- RIGIRVSP igtfq ...
	... k P-LGVKLPP yf--dlvhfdimaeilnqfpltyvsnv-nsig----nglfidpeaesv ...
	... yl-1nqfldphsntrtdeyggssienrarftlevvdalveaigh- KVGLRLSP ygvfn ...
	... yl-plqflnpyynkrtdkyggslenrarfwletlekvhavgsdc AIATRF -- GV dt ...
	... kv PLYVKLSP nv-tdivpiakaveaagadgltmintl-----mgvrfdlktrqp ...
default	... gsvenrarlvlevvdavcnewsad- RIGIRVSP igtfqnvdngpnee--adalyl--- ...
	... ydfeatekllke-----vftfftk- PLGVKLPP yf-----dlvhfdim ...
	... gsienrarftlevvdalveaigh- KVGLRLSP ygvfnmsggaetgivaqyayvage ...
	... gslenrarfwletlekvhavgsdc AIATRGV -----dtvygpgq ...
	... tdpevaaalvka-----ckavskv- PLYVKLSP nvt-----divpiaka ...
alternate	... gsvenrarlvlevvdavcnewsad- RIGIRVSP igtfqnvdngpnee--adalyl--- ...
	... ydfeatekllke-----vftfftk- PLGVKLPP yf-----dlvhfdim ...
	... gsienrarftlevvdalveaigh- KVGLRLSP ygvfnmsggaetgivaqyayvage ...
	... gslenrarfwletlekvhavgsdc AIATRGV -----dtvygpgq ...
	... tdpevaaalvka-----ckavskv- PLYVKLSP nvt-----divpiaka ...

Motivation

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

Motivation

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

- Alignment parameters **model** mutation rates

Motivation

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

- Alignment parameters **model** mutation rates
- A **single choice** of parameters may not be best

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

Motivation

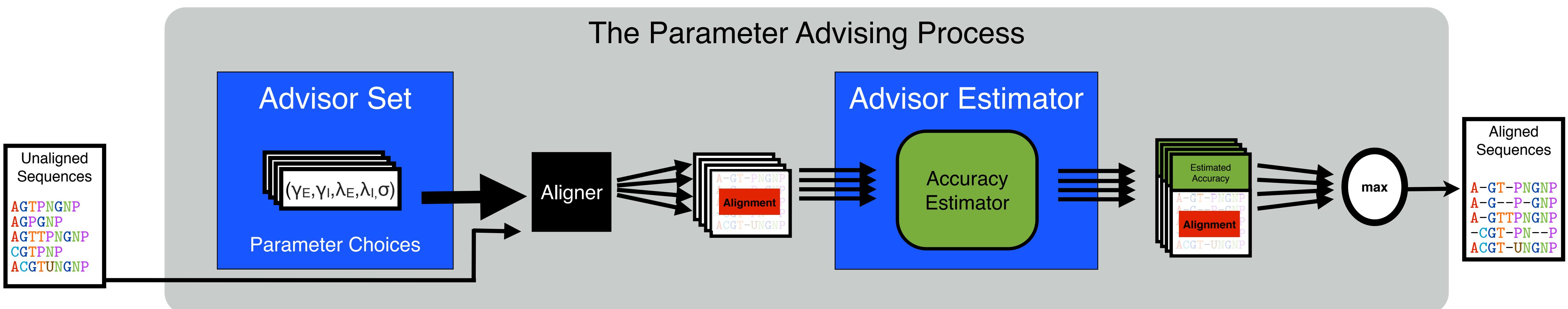
Proteins can have **different mutation rates** along their length

- 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

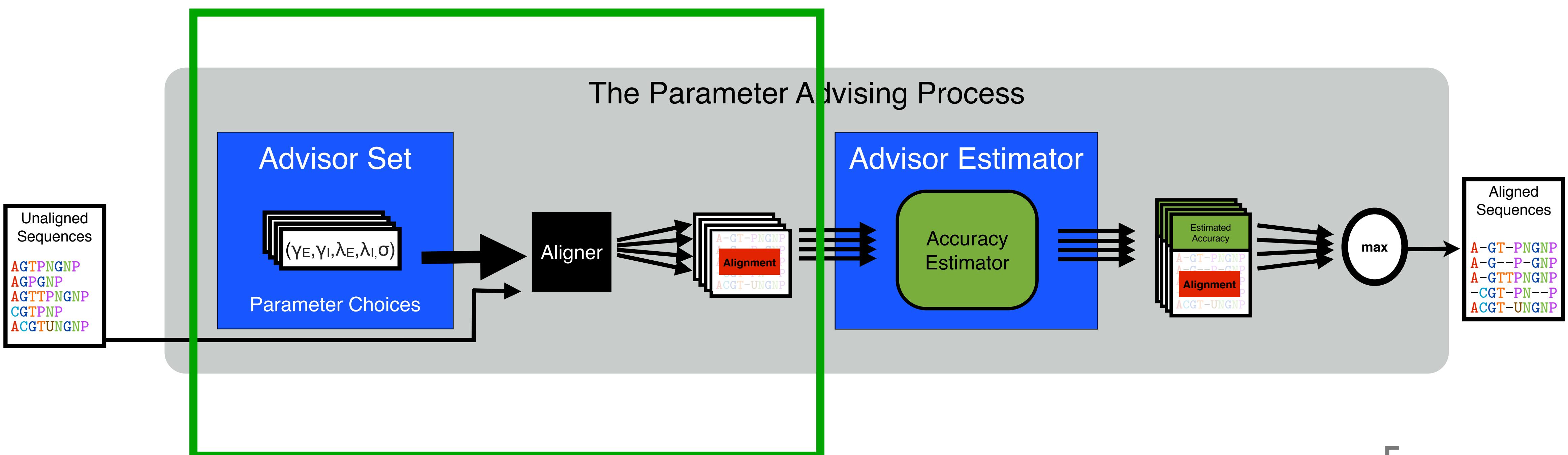
Advising for unaligned input sequences



Parameter advising

Advising for unaligned input sequences

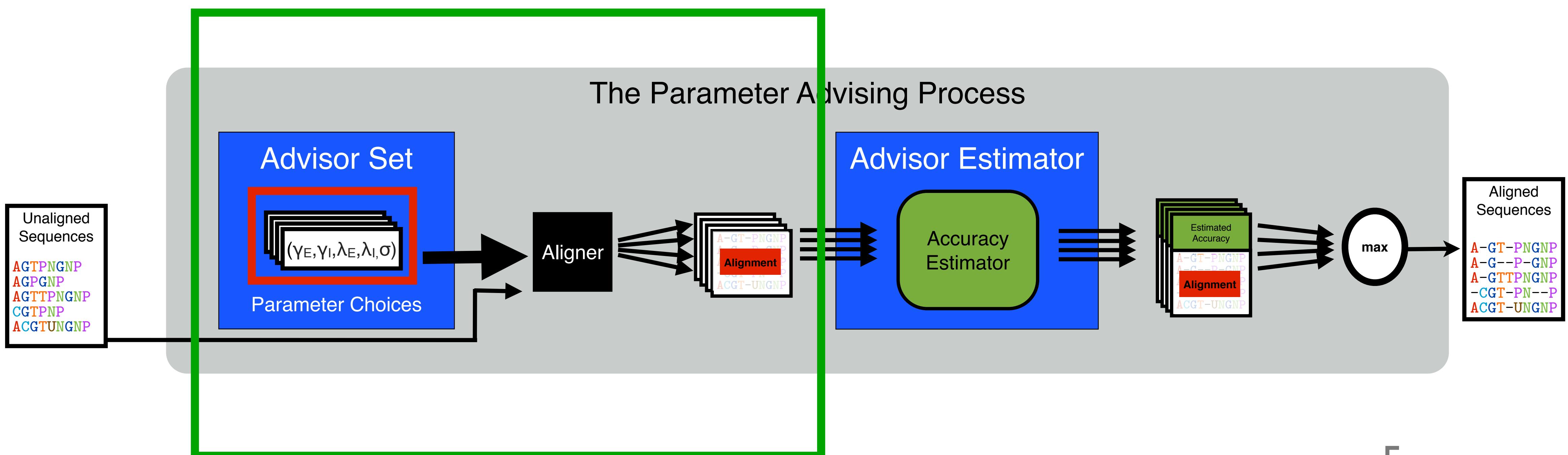
- aligns sequences using each parameter choice from a set,



Parameter advising

Advising for unaligned input sequences

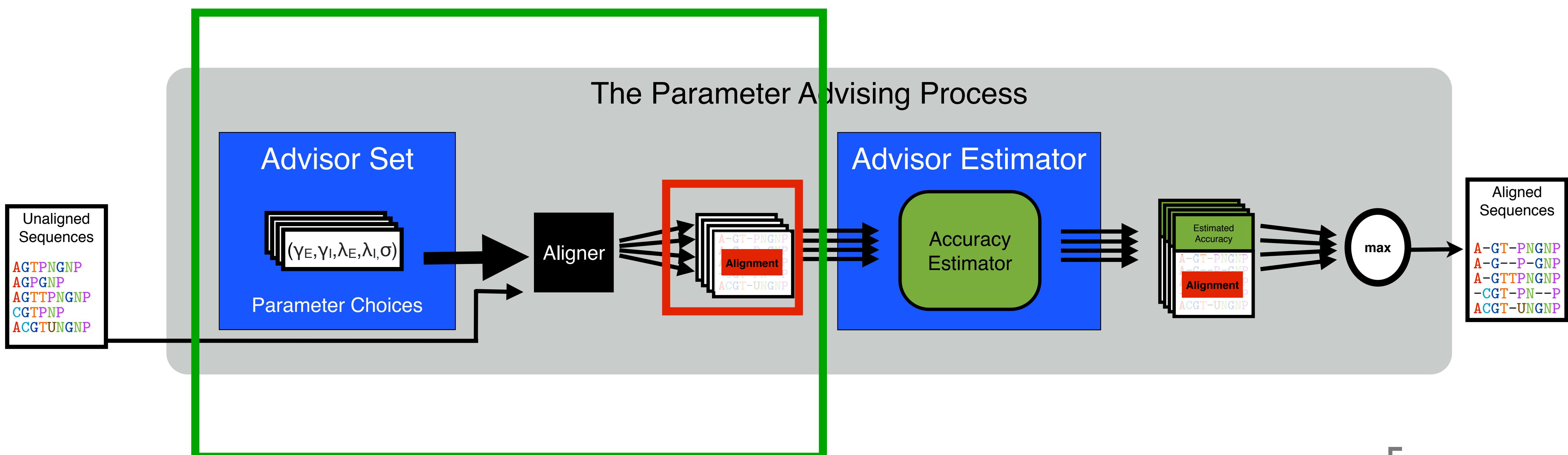
- aligns sequences using each parameter choice from a set,



Parameter advising

Advising for unaligned input sequences

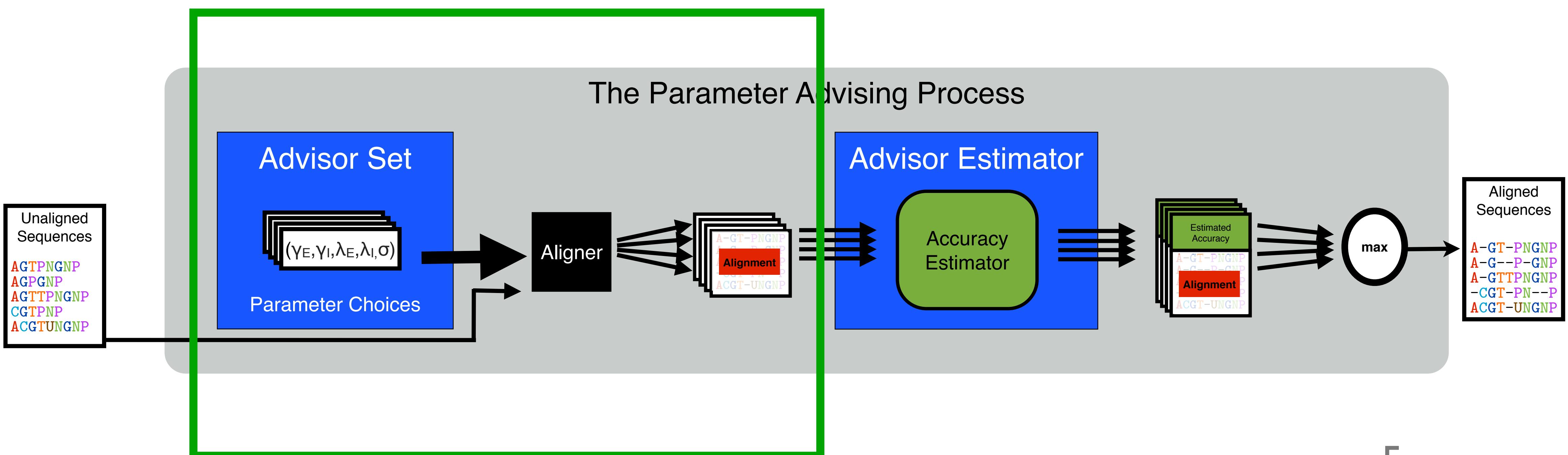
- aligns sequences using each parameter choice from a set,



Parameter advising

Advising for unaligned input sequences

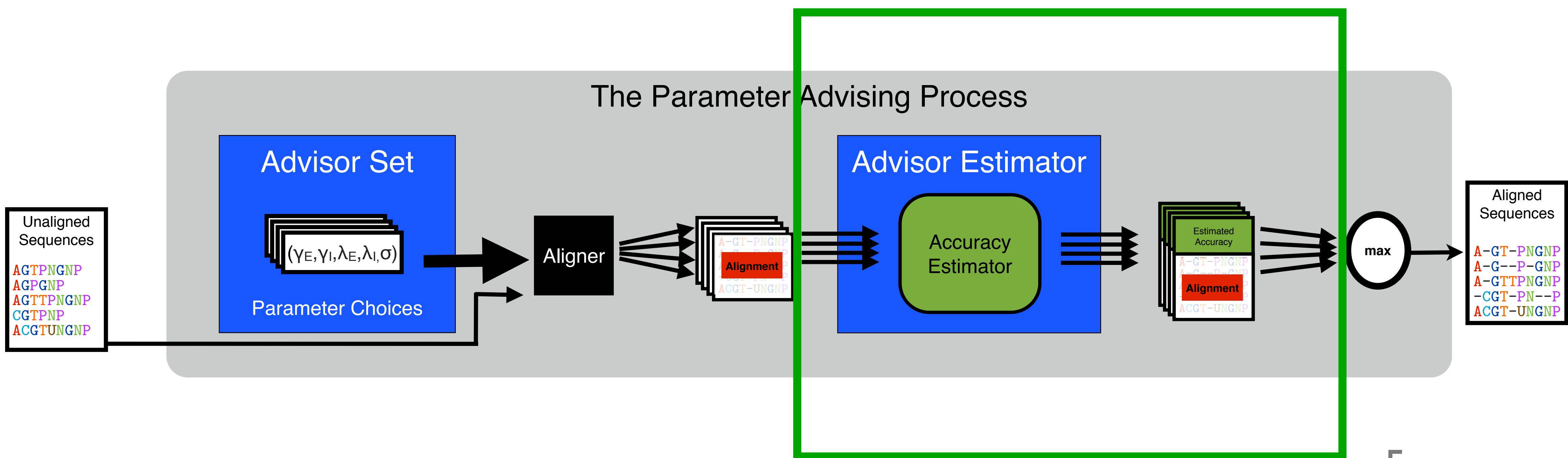
- aligns sequences using each parameter choice from a set,



Parameter advising

Advising for unaligned input sequences

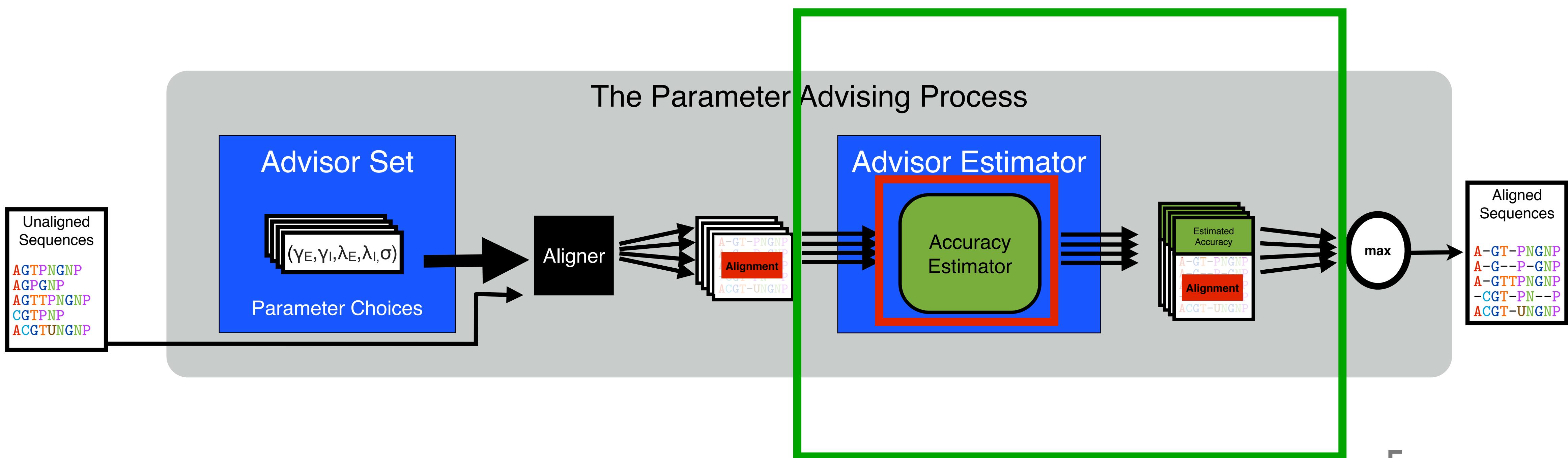
- aligns sequences using each parameter choice from a set,
- assigns an estimated accuracy to each alignment, and



Parameter advising

Advising for unaligned input sequences

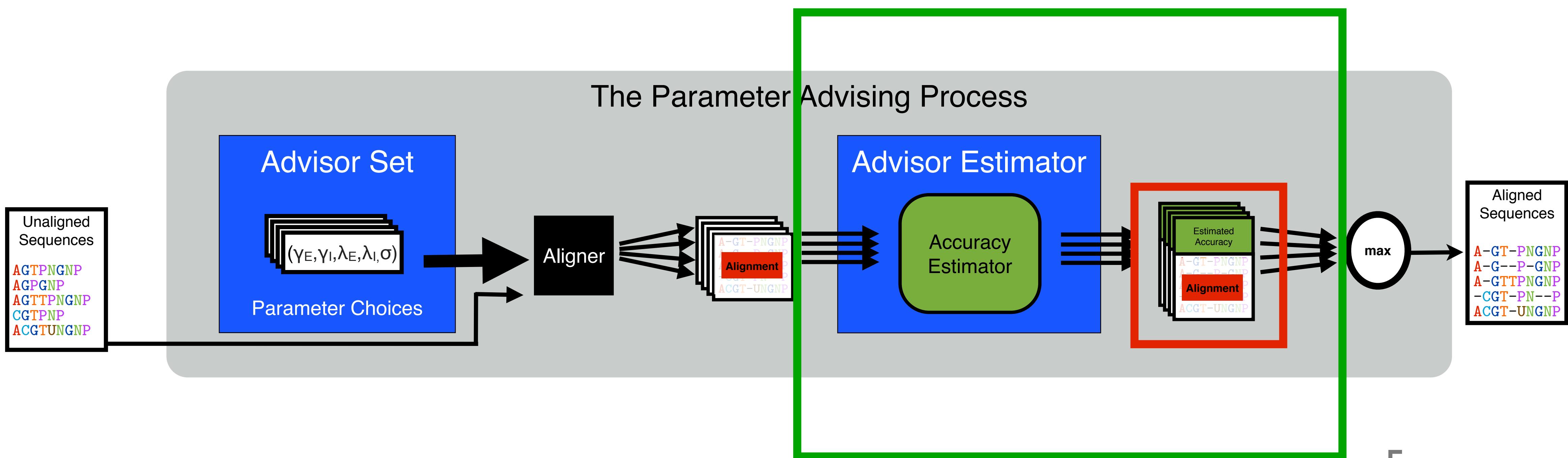
- aligns sequences using each parameter choice from a set,
- assigns an estimated accuracy to each alignment, and



Parameter advising

Advising for unaligned input sequences

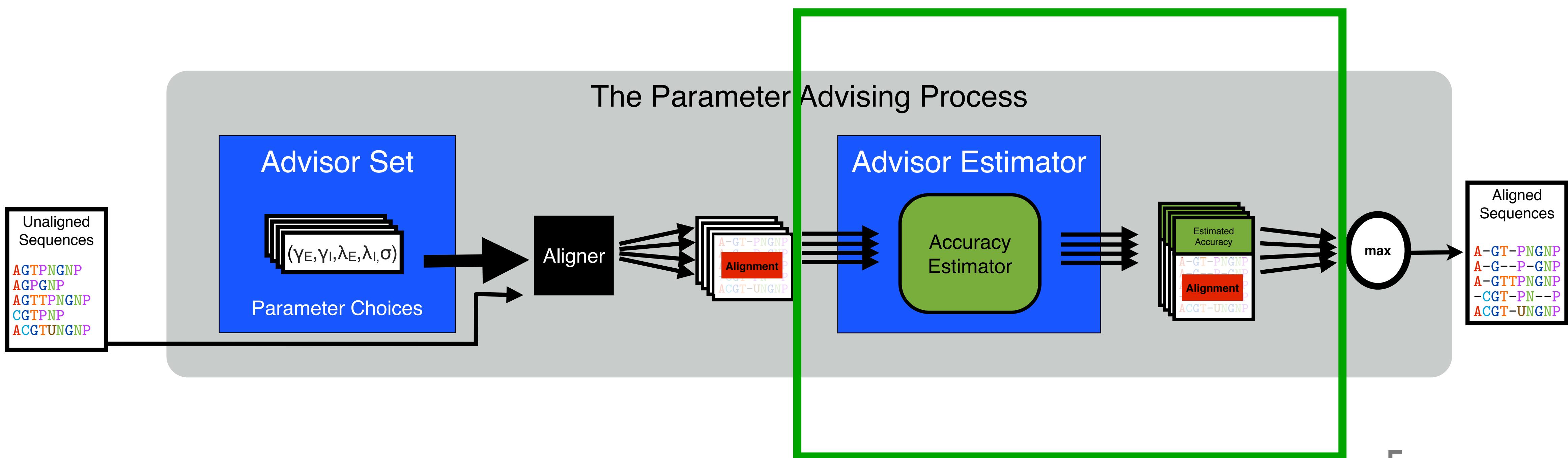
- aligns sequences using each parameter choice from a set,
- assigns an estimated accuracy to each alignment, and



Parameter advising

Advising for unaligned input sequences

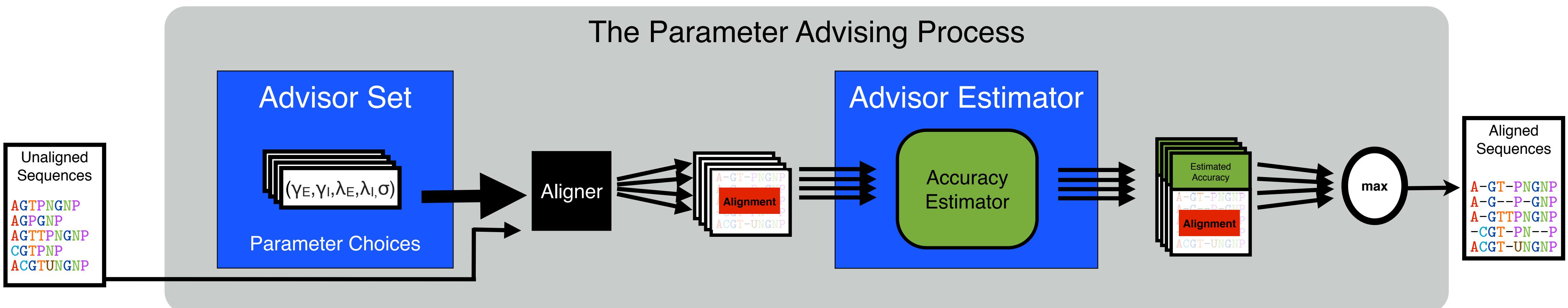
- aligns sequences using each parameter choice from a set,
- assigns an estimated accuracy to each alignment, and



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.

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 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 aligned residue pairs 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 aligned residue pairs 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 aligned residue pairs 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 aligned residue pairs 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 aligned residue pairs 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 ...	66%
... a S H l t a S - H l t ...	Accuracy

↑ ↑

- accuracy is the fraction of aligned residue pairs from the reference that are in the computed alignment,
- measured on the core columns of the reference.

Accuracy estimation

Our estimator **Facet** (“Feature-based **AC**curacy **E**s**T**imator”)

Accuracy estimation

Our estimator **Facet** (“Feature-based A^Ccuracy E^Stimator”)

- estimates accuracy by a **polynomial** on feature functions,

Accuracy estimation

Our estimator **Facet** (“Feature-based A^Ccuracy EsTimator”)

- estimates accuracy by a **polynomial** on feature functions,
- uses **novel features** that are efficient to evaluate,

Accuracy estimation

Our estimator **Facet** (“Feature-based A^Ccuracy EsTimator”)

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

Advisor sets

An advisor set should

Advisor sets

An **advisor set** should

- consist of **complementary** parameter settings,

Advisor sets

An **advisor set** should

- consist of **complementary** parameter settings,
- produce at least one **good alignment** for every input,

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**.

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

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

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**,

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**,

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**.

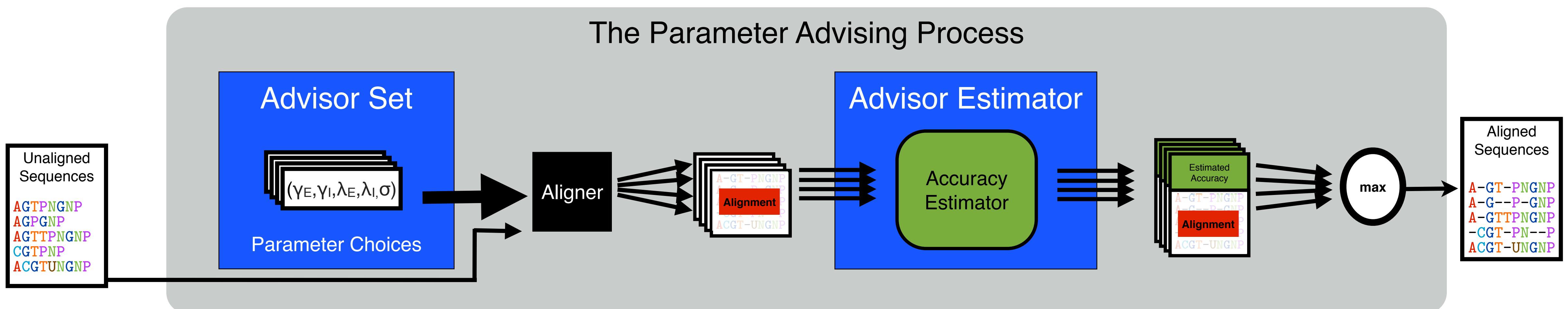
An **oracle set** is

- 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.

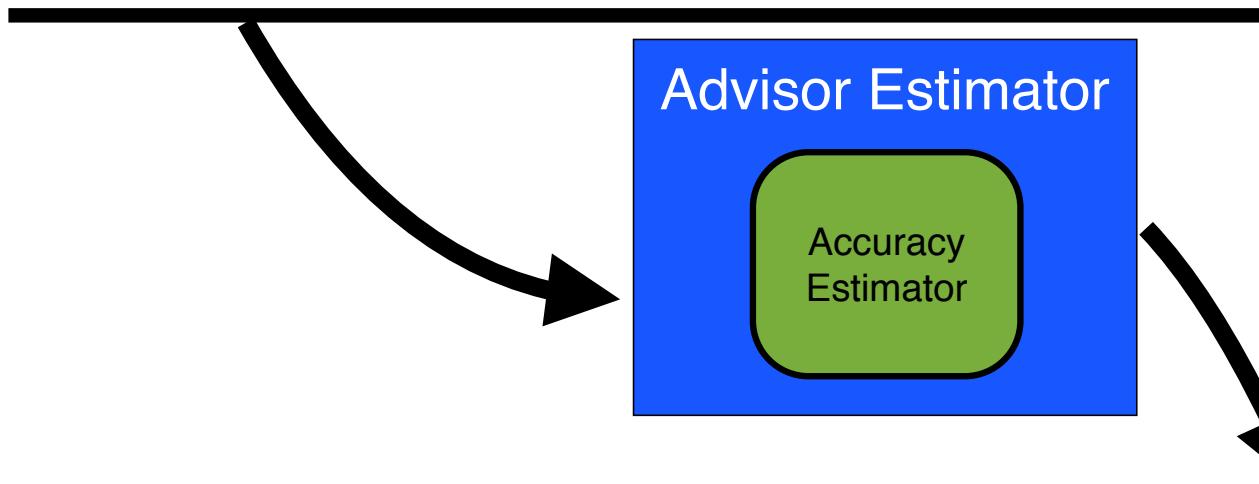


Adaptive local realignment

rkeyag**LYHEVAQAHGVDSQV**rq**MKFGLFFLFDTLAV**yenhsnngvldqmsegrfafhkiindafttgychpnnd**PLVF**rwddsnaq**HKLILLVNQNVGEAARAEARV**yleefvresysnt
kkaqlid**LYNEVATEHGYDVTKI**d-**MKFGNFLLFDTVWL**lehhftefgllldqmskgrfrfydlmkegfnegyiaadne**PMIL**swiinthe**HCLSYITSVDHDSNRACKDICRN**flghwy-dsyvna
rielln**HYQAAAAKFNVDIANV**r-----mtk**WNYGVFFLYDVVA--F**sehhidksyn**PLLF**kwddsqqk**HRLMLFVNNDNPQAKAELSI**yledyl--sytqa
rlklls**FYNASASKYNKNIDLV**r-----mnk**WNYGVFFVYDVIN--I**ddhylvkkds**PLVF**kwddinee**HQLMLHVNVNEAETVAKEELKL**yienyv--actqp

Adaptive local realignment

rkeyag**LYHEVAQAHGVDSQV**rq**MKFGLFFLFDTLAV**yenhsnngvldqmsegrfafhkiindafttgcypnnd**PLVF**rwddsnaq**HKLILLVNQNVGEAARAEARV**yleefvresysnt
kkaqlid**LYNEVATEHGYDVTKI**d-**MKFGNFLLFDTVWL**lehhftefgllldqmskgrfrfydlmkefnegyiaadne**PMIL**swiinthe**HCLSYITSVDHDSNRACKDICRN**flghwy-dsyvna
rielln**HYQAAAAKFNVDIANV**r-----mtk**WNYGVFFLYDVVA--F**sehhidksyn**PLLF**kwddsqqk**HRLMLFVNNDNPQAKAEISI**yledyl--sytqa
rlklls**FYNASASKYNKNIDLV**r-----mnk**WNYGVFFVYDVIN--I**ddhylvkkds**PLVF**kwddinee**HQLMLHVNVNEAETVAKEELKL**yienyv--actqp



Adaptive local realignment



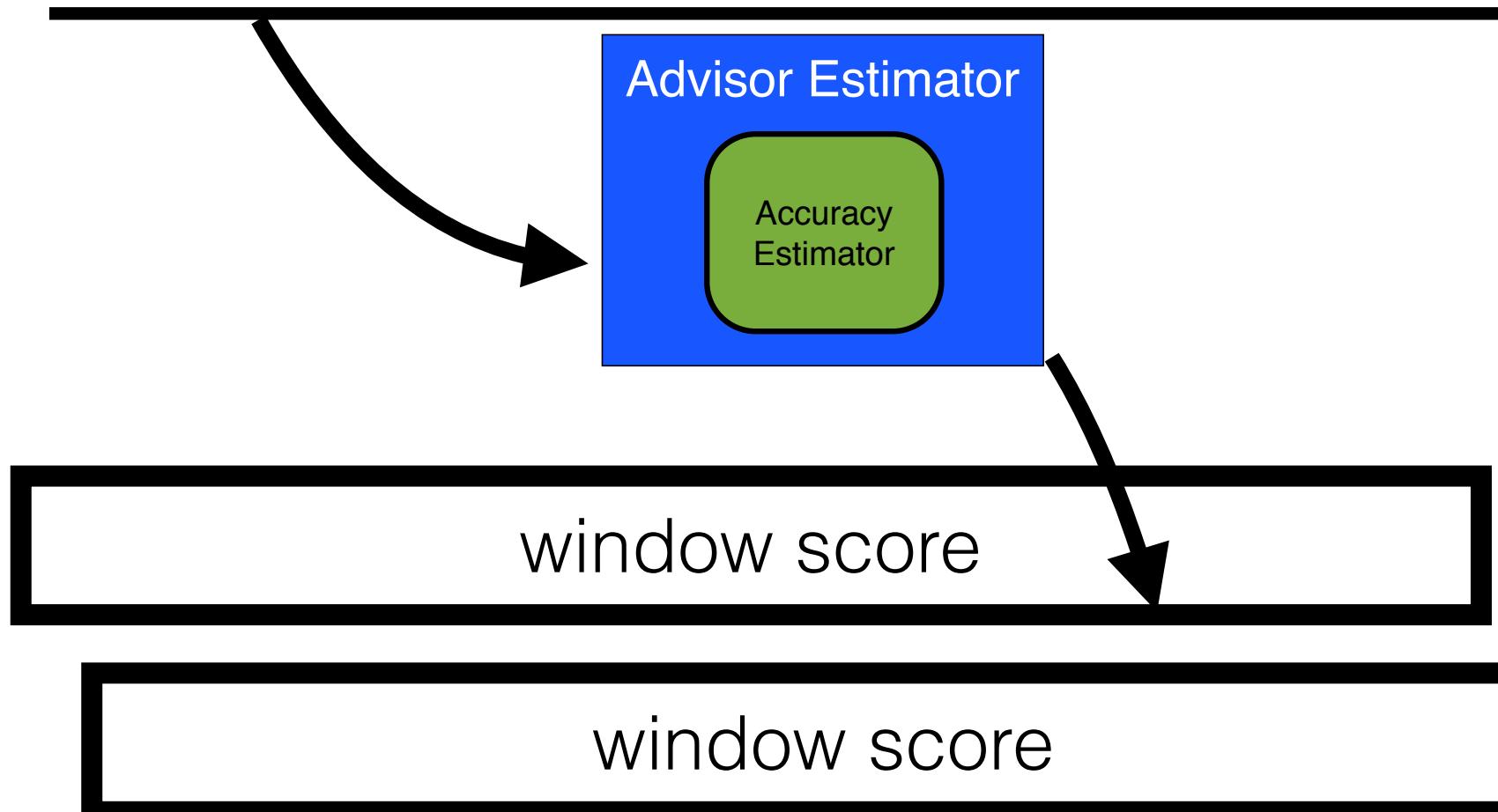
Adaptive local realignment

rkeyag**LYHEVAQAHGVDSQV**rq**MKFGLFFLFDTLAV**yenhsnngvldqmsegrfafhkiindafttgcchpnnd**PLVF**rwddsnaq**HKLILLVNQNVGEAARAEARV**yleefvresysnt
kkaqlid**LYNEVATEHGYDVTKI**d-**MKFGNFLLFDTVWL**lehhftefgllldqmskgrfrfydlmkegfnegyiaadne**PMIL**swiinthe**HCLSYITSVDHDSNRACKDICRN**flghwy-dsyvna
rielln**HYQAAAAKFNVDIANV**r-----mtk**WNYGVFFLYDVVA--F**sehhidksyn**PLLF**kwddsqqk**HRLMLFVNNDNPTQAKAELSI**yledyl--sytqa
rlklls**FYNASASKYNKNIDLV**r-----mnk**WNYGVFFVYDVIN--I**ddhylvkkds**PLVF**kwddinee**HQLMLHVNVNEAETVAKEELKL**yienyv--actqp

window score

Adaptive local realignment

rkeyag**LYHEVAQAHGVDSQV**rq**MKFGLFFLFDTLAV**yenhsnngvldqmsegrfafhkiindafttgcypnnd**PLVF**rwddsnaq**HKLILLVNQNVGEAARAEARV**yleefvresysnt
kkaqld**LYNEVATEHGYDVTKI**d-**MKFGNFLLFDTVWL**lehhftefgllldqmskgrfrfydlmkefnegyiaadne**PMIL**swiinthe**HCLSYITSVDHDSNRACKDICRN**flghwy-dsyvna
rielln**HYQAAAAKFNVDIANV**r-----mtk**WNYGVFFLYDVVA--F**sehhidksyn**PLLF**kwddsqqk**HRLMLFVNNDNPQAKAELSI**yledyl--sytqa
rlklls**FYNASASKYNKNIDLV**r-----mnk**WNYGVFFVYDVIN--I**ddhylvkkds**PLVF**kwddinee**HQLMLHVNVNEAETVAKEELKL**yienyv--actqp



Adaptive local realignment

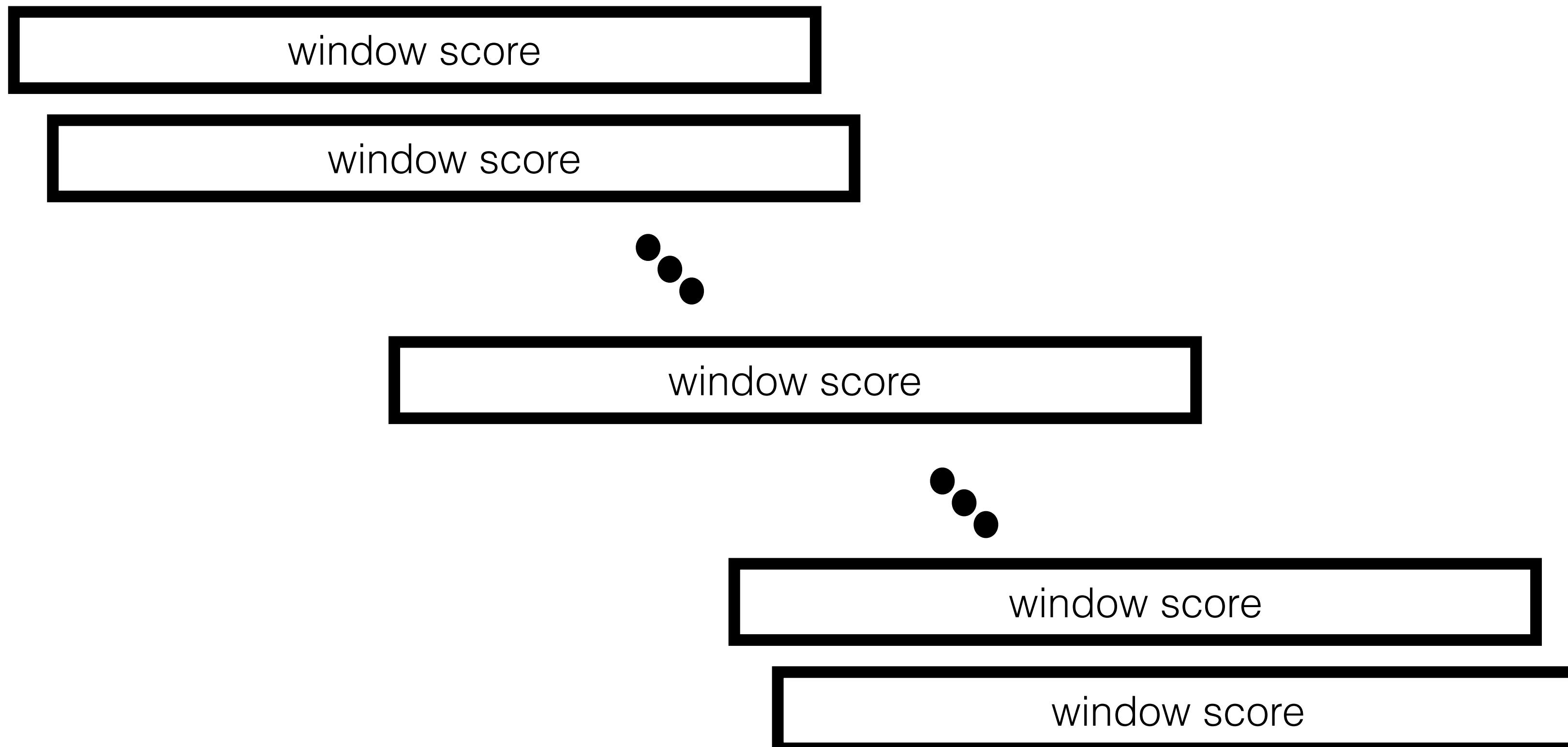
rkeyag**LYHEVAQAHGVDSQV**rq**MKFGLFFLFDTLAV**yenhsnngvldqmsegrfafhkiindafttgcchpnnd**PLVF**rwddsnaq**HKLILLVNQNVGEAARAEARV**yleefvresysnt
kkaqlid**LYNEVATEHGYDVTKI**d-**MKFGNFLLFDTVWL**lehhftefgllldqmskgrfrfydlmkegfnegyiaadne**PMIL**swiinthe**HCLSYITSVDHDSNRACKDICRN**flghwy-dsyvna
rielln**HYQAAAAKFNVDIANV**r-----mtk**WNYGVFFLYDVVA--F**sehhidksyn**PLLF**kwddsqk**HRLMLFVNNDNPTQAKAELSI**yledyl--sytqa
rlklls**FYNASASKYNKNIDLV**r-----mnk**WNYGVFFVYDVIN--I**ddhylvkkds**PLVF**kwddinee**HQLMLHVNVNEAETVAKEELKL**yienyv--actqp

window score

window score

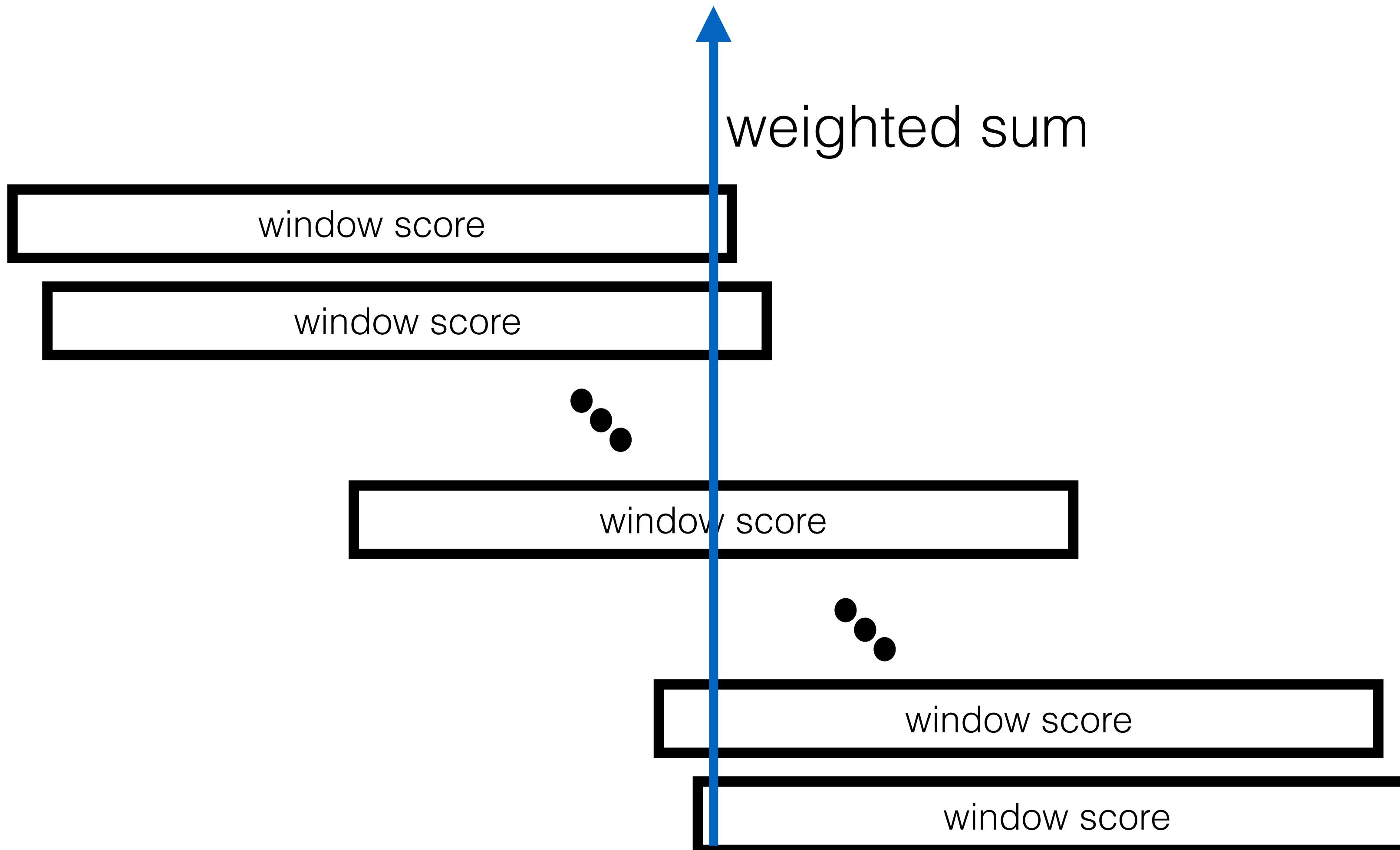
Adaptive local realignment

rkeyag**LYHEVAQAHGVDSQV**rq**MKFGLFFLFDTLAV**yenhsnngvldqmsegrfafhkiindafttgcchpnnd**PLVF**rwddsnaq**HKLILLVNQNVGEAARAEARV**yleefvresysnt
kkaqld**LYNEVATEHGYDVTKI**d-**MKFGNFLLFDTVWL**lehhftefgllldqmskgrfrfydlmkegfnegyiaadne**PMIL**swiinthe**HCLSYITSVDHDSNRAKDICRN**flghwy-dsyvna
rielln**HYQAAAAKFNVDIANV**r-----mtk**WNYGVFFLYDVVA--F**sehhidksyn**PLLF**kwddsqk**HRLMLFVNNDNPTQAKAELSI**yledyl--sytqa
rlklls**FYNASASKYNKNIDLV**r-----mnk**WNYGVFFVYDVIN--I**ddhylvkkds**PLVF**kwddinee**HQLMLHVNVNEAETVAKEELKL**yienyv--actqp



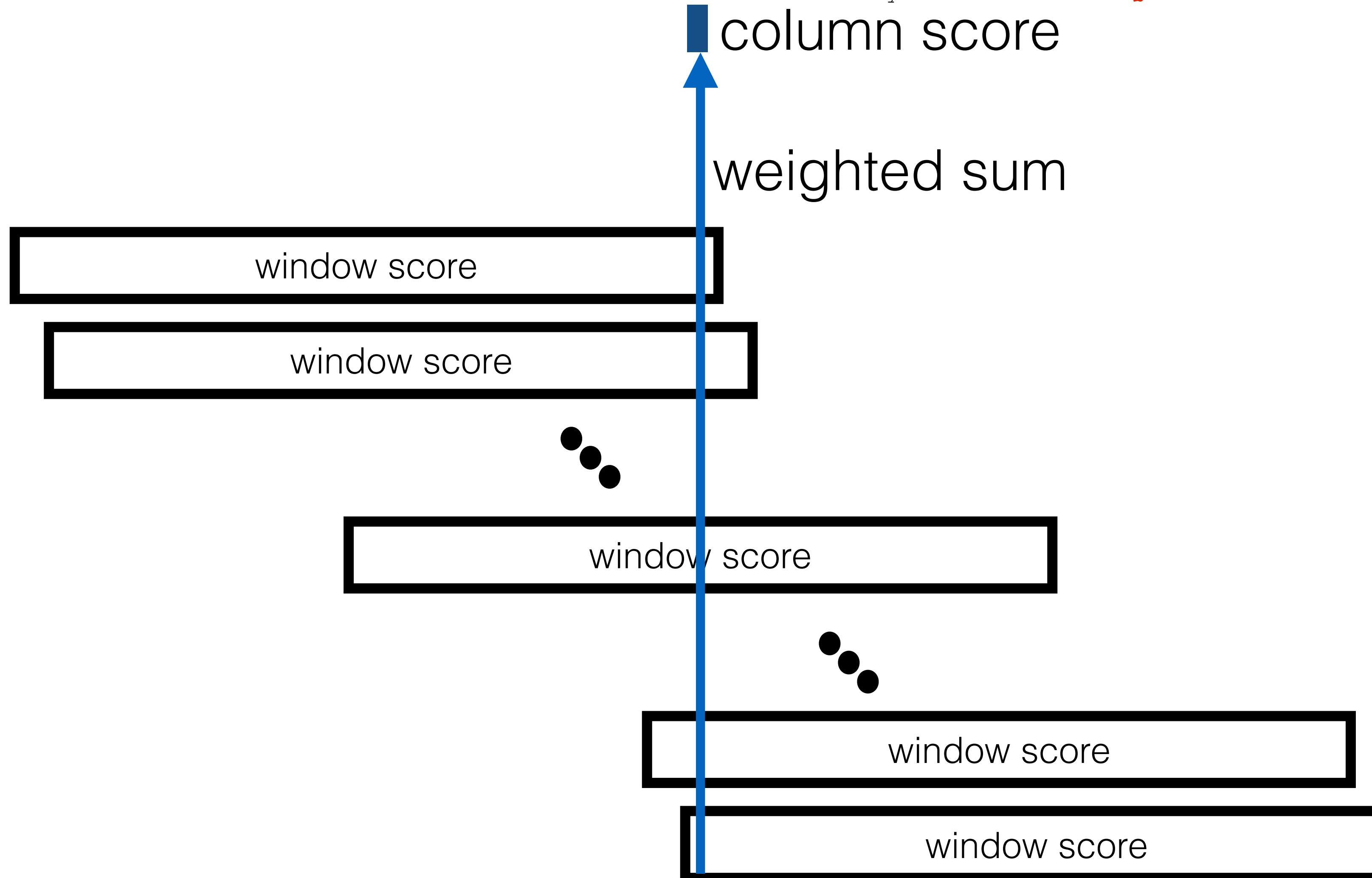
Adaptive local realignment

rkeyag**LYHEVAQAHGVDSQV**rq**MKFGLFFLFDTLAV**yenhsnngvldqmsegrfafhkiindafttgcypnnd**PLVF**rwddsnaq**HKLILLVNQNVGEAARAEARV**yleefvresysnt
kkaqld**LYNEVATEHGYDVTKI**d-**MKFGNFLLFDTVWL**lehhftefgllldqmskgrfrfydlmkefnegyiaadne**PMIL**swiinthe**HCLSYITSVDHDSNRAKDICRN**flghwy-dsyvna
rielln**HYQAAAAKFNVDIANV**r-----mtk**WNYGVFFLYDVVA--F**sehhidksyn**PLLF**kwddsqk**HRLMLFVNNDNPTQAKAELSI**yledyl--sytqa
rlklls**FYNASASKYNKNIDLV**r-----mnk**WNYGVFFVYDVIN--I**ddhylvkkds**PLVF**kwddinee**HQLMLHVNVNEAETVAKEELKL**yienyv--actqp



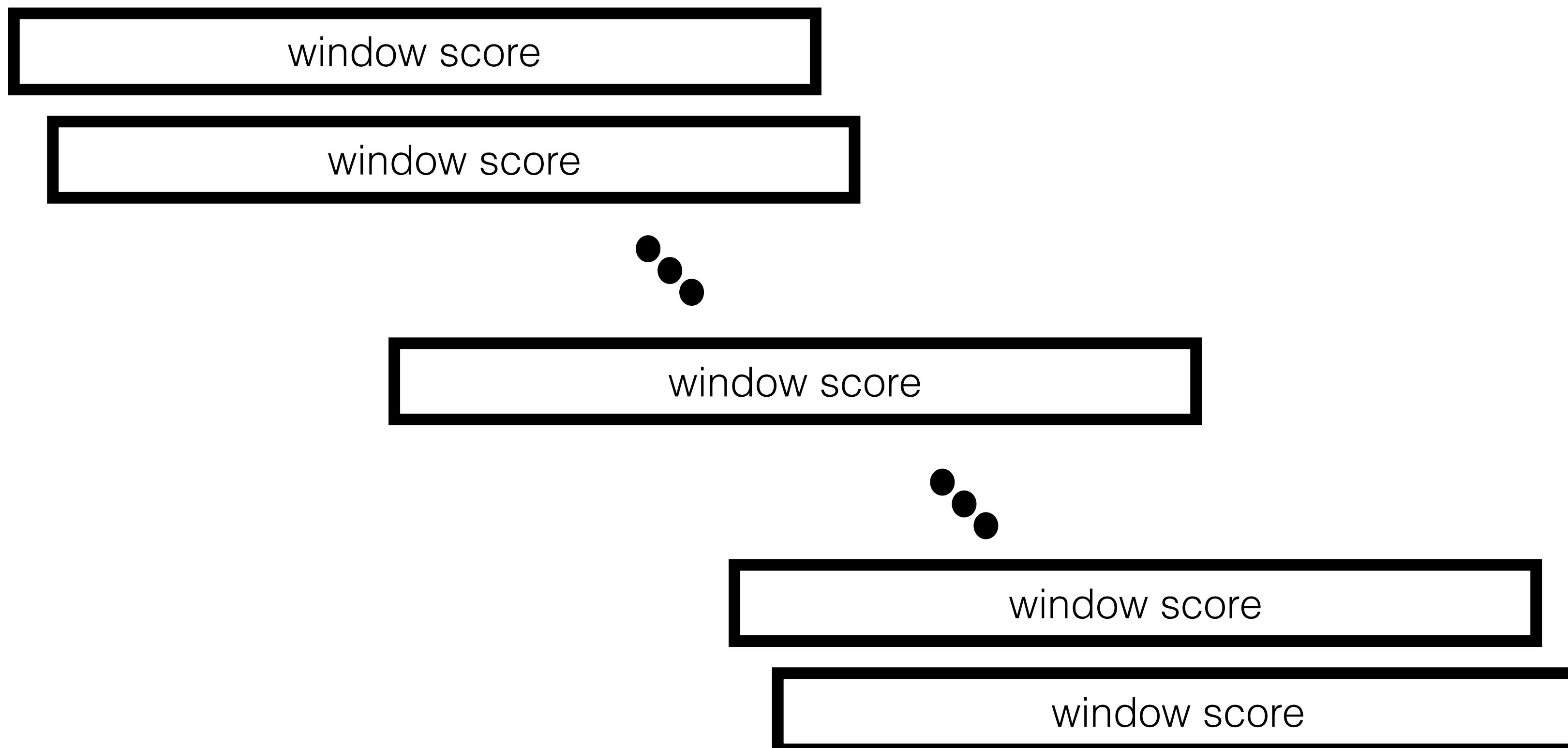
Adaptive local realignment

rkeyag**LYHEVAQAHGVDSQV**rq**MKFGLFFLFDTLAV**yenhsnngvvlldqmsegrfafhkiindafttgcypnnd**PLVF**rwddsnaq**HKLILLVNQNVGEAARAEARV**yleefvresysnt
kkaqld**LYNEVATEHGYDVTKI**d-**MKFGNFLLFDTVWL**lehhftefgllldqmskgrfrfydlmkegfnegyiaadne**PMIL**swiinthe**HCLSYITSVDHDSNRAKDICRN**flghwy-dsyvna
rielln**HYQAAAAKFNVDIANV**r-----mtk**WNYGVFFFLYDVVA--F**sehhidksyn**PLLF**kwddsqqk**HRLMLFVNNDNPQAKAELSI**yledyl--sytqa
rlklls**FYNASASKYNKNIDLV**r-----mnk**WNYGVFFVYDVIN--I**ddhylvkkds**PLVF**kwddinee**HQLMLHVNVNEAETVAKEELKL**yienyv--actqp



Adaptive local realignment

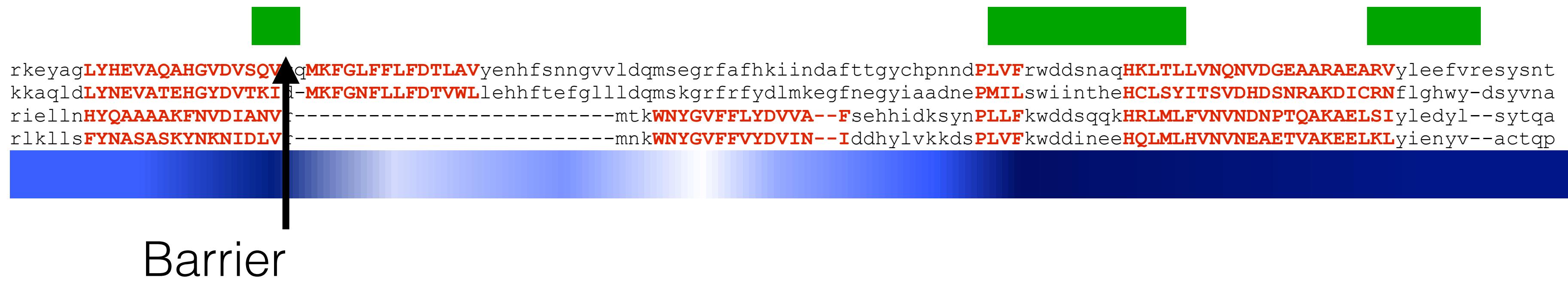
rkeyag**LYHEVAQAHGVDSQV**rq**MKFGLFFLFDTLAV**yenhsnngvldqmsegrfafhkiindafttgcypnnd**PLVF**rwddsnaq**HKLILLVNQNVDGEAARAEARV**yleefvresysnt
kkaql*d***LYNEVATEHGYDVTKI***d*-**MKFGNFLLFDTVWL**lehhftefgllldqmskgrfrfydlmkefnegyiaadne**PMIL**swiinthe**HCLSYITSVDHDSNRAKDICRN**flghwy-dsyvna
rielln**HYQAAAAKFNVDIANV**r-----mtk**WNYGVFFLYDVVA--F**sehhidksyn**PLLF**kwddsqqk**HRLMLFVNNDNPTQAKAELSI**yledyl--sytqa
rlklls**FYNASASKYNKNIDLV**r-----mnk**WNYGVFFVYDVIN--I**ddhylvkkds**PLVF**kwddinee**HQLMLHVNVNEAETVAKEELKL**yienyv--actqp



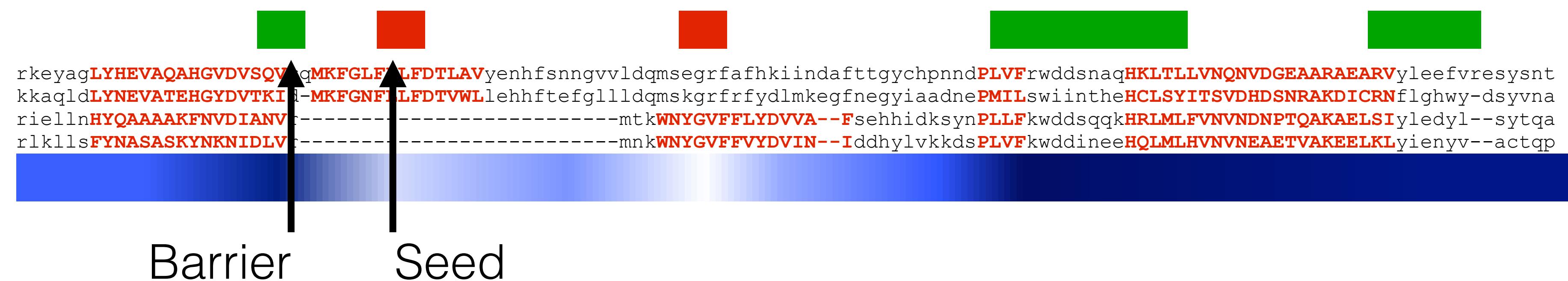
Adaptive local realignment

rkeyag**LYHEVAQAHGVDSQV**rq**MKFGLFFLFDTLAV**yenhsnngvldqmsegrfafhkiindafttgcchpnnd**PLVF**rwddsnaq**HKLILLVNQNVGEAARAEARV**yleefvresysnt
kkaqlid**LYNEVATEHGYDVTKI**d-**MKFGNFLLFDTVWL**lehhftefgllldqmskgrfrfydlmkegfnegyiaadne**PMIL**swiinthe**HCLSYITSVDHDSNRACKDICRN**flghwy-dsyvna
rielln**HYQAAAAKFNVDIANV**r-----mtk**WNYGVFFLYDVVA--F**sehhidksyn**PLLF**kwddsqqk**HRLMLFVNNDNPTQAKAELSI**yledyl--sytqa
rlklls**FYNASASKYKNIDLV**r-----mnk**WNYGVFFVYDVIN--I**ddhylvkkds**PLVF**kwddinee**HQLMLHVNVNEAETVAKEELKL**yienyv--actqp

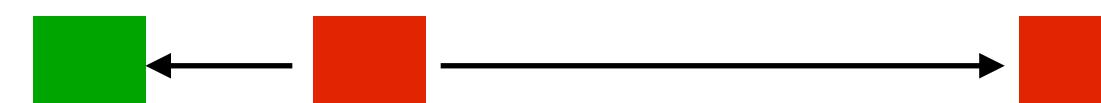
Adaptive local realignment



Adaptive local realignment



Adaptive local realignment



rkeyag**LYHEVAQAHGVDSQV**rq**MKFGLFFLFDTLAV**yenhsnngvldqmsegrfafhkiindafttgychpnnd**PLVF**rwddsnaq**HKL****TLLVNQNVDGEAARAEARV**yleefvresysnt
kkaqld**LYNEVATEHGYDVTKI**d-**MKFGNFLLFDTVWL**lehhftefgllldqmskgrfrfydlmkegfnegyiaadne**PMIL**swiinthe**HCLSYITSVDHDSNRACKDICRN**flghwy-dsyvna
rielln**HYQAAAAKFNVDIANV**r-----mtk**WNYGVFFLYDVVA--F**sehhidksyn**PLLF**kwddsqqk**HRLMLFVNNDNPQAKAELSI**yledyl--sytqa
rlklls**FYNASASKYNKNIDLV**r-----mnk**WNYGVFFVYDVIN--I**ddhylvkkds**PLVF**kwddinee**HQLMLHVNVNEAETVAKEELKL**yienyv--actqp



Adaptive local realignment

rkeyag**LYHEVAQAHGVDSQV**rq**MKFGLFFLFDTLAV**yenhsnngvldqmsegrfafhkiindafttgychpnnd**PLVF**rwddsnaq**HKL****TLLVNQNVDGEAARAEARV**yleefvresysnt
kkaqlid**LYNEVATEHGYDVTKI**d-**MKFGNFLLFDTVWL**lehhftefgllldqmskgrfrfydlmkegfnegyiaadne**PMIL**swiinthe**HCLSYITSVDHDSNRAKDICRN**flghwy-dsyvna
rielln**HYQAAAAKFNVDIANV**r-----mtk**WNYGVFFLYDVVA--F**sehhidksyn**PLLF**kwddsqqk**HRLMLFVNNDNPQAKAELSI**yledyl--sytqa
rlklls**FYNASASKYNKNIDLV**r-----mnk**WNYGVFFVYDVIN--I**ddhylvkkds**PLVF**kwddinee**HQLMLHVNVNEAETVAKEELKL**yienyv--actqp

Adaptive local realignment



Adaptive local realignment

rkeyag**LYHEVAQAHGVDSQV**rq**MKFGLFFLFDTLAV**yenhsnngvldqmsegrfafhkiindafttgcypnnd**PLVF**rwddsnaq**HKL****TLLVNQNVDGEAARAEARV**yleefvresysnt
kkaqlid**LYNEVATEHGYDVTKI**d-**MKFGNFLLFDTVWL**lehhftefgllldqmskgrfrfydlmkegfneyiaadne**PMIL**swiinthe**HCLSYITSVDHDSNRAKDICRN**flghwy-dsyvna
rielln**HYQAAAAKFNVDIANV**r-----mtk**WNYGVFFLYDVVA--F**sehhidksyn**PLLF**kwddsqqk**HRLMLFVNNDNPQAKAELSI**yledyl--sytqa
rlklls**FYNASASKYNKNIDLV**r-----mnk**WNYGVFFVYDVIN--I**ddhylvkkds**PLVF**kwddinee**HQLMLHVNVNEAETVAKEELKL**yienyv--actqp

Adaptive local realignment

rkeyag**LYHEVAQAHGVDSQV**rq**MKFGLFFLFDTLAV**yenhsnngvvldqmsegrfafhkiindafttgychpnnd**ELVF**rwddsnaq**HKL****TLLVNQNVDGEAARAEARV**yleefvresysnt
kkaqlid**LYNEVATEHGYDVTKI**d**-MKFGNFLLFDTVWL**lehhftefgllldqmskgrfrfydlmkegfnegyiaadne**PMIL**swiinthe**HCLSYITSVDHDSNRAKDICRN**flghwy-dsyvna
rielln**HYQAAAAKFNVDIANV**r-----mtk**WNYGVFFLYDVVA--F**sehhidksyn**ELLF**kwddsqqk**HRLMLFVNVNDNPTQAKAELSI**yledyl--sytqa
rlklls**FYNASASKYNKNIDLV**r-----mnk**WNYGVFFVYDVIN--I**ddhylvkkds**ELVF**kwddineee**HQLMLHVNVNEAETVAKEELKL**yienyv--actqp

Adaptive local realignment

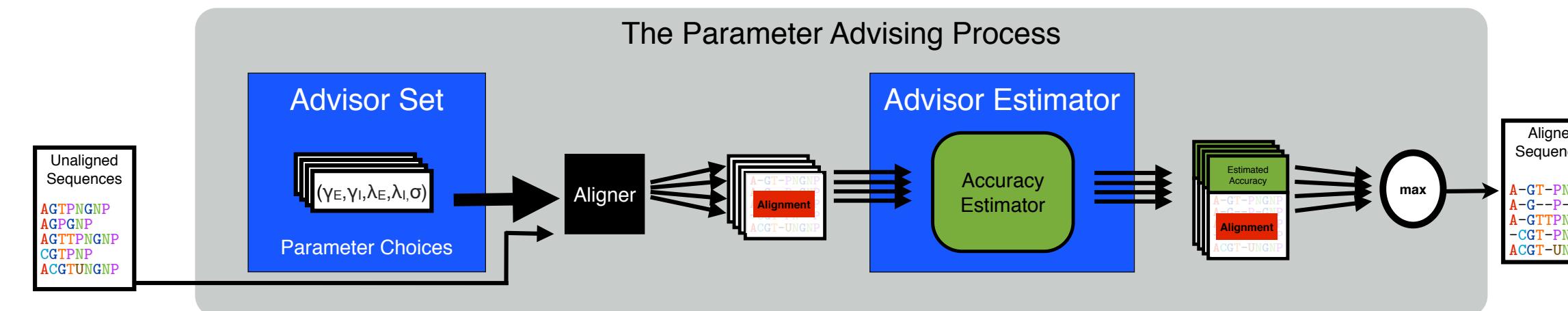
rkeyag **LYHEVAQAHGVDSQV** r q **MKFGLFFLFDTLAV** yenhfsnngvvldqmsegrfafhkiindafttgychpnnd **ELVF** rwddsnaq **HKLTLLVNQNVGEAARAEARV** yleefvresysnt
kkaql d **LYNEVATEHGYDVTKI** d - **MKFGNFLLFDTVWL** lehhftefgllldqmskgrfrfydlmkegfnegyiaadne **PMIL** swi in the **HCLSYITSVDHDSNRACKDICRN** flghwy-dsyvna
rielln **HYQAAAAKFNVDIANV** r -----mtk **WNYGVFFFLYDVVA--F** sehhidksyn **ELLF** kwddsqqk **HRLMLFVNNDNPTQAKAELSI** yledyl--sytqa
rlklls **FYNASASKYNKNIDLV** r -----mnk **WNYGVFFVYDVIN--I** ddhylvkkds **ELVF** kwddineee **HQLMLHVNVNEAETVAKEELKL** ienyv--actqp

q **MKFGLFFLFDTLAV** yenhfsnngvvldqmsegrfafhkiindafttgychpnnd
- **MKFGNFLLFDTVWL** lehhftefgllldqmskgrfrfydlmkegfnegyiaadne
-----mtk **WNYGVFFFLYDVVA--F** sehhidksyn
-----mnk **WNYGVFFVYDVIN--I** ddhylvkkds

Adaptive local realignment

rkeyag **LYHEVAQAHGVDSQV** q **MKFGLFFLFDTLAV** yenhfsnngvvldqmsegrfafhkiindafttgcypnnd **ELVF** rwddsnaq **HKLTLLVNQNVGEAARAEARV** yleefvresysnt
kkaqlid **LYNEVATEHGYDVTKI** d - **MKFGNFLLFDTVWL** lehhftefgllldqmskgrfrfydlmkegfnegyiaadne **PMIL** swiinthe **HCLSYITSVDHDSNRAKDICRN** flghwy-dsyvna
rielln **HYQAAAAKFNVDIANV** r -----mtk **WNYGVFFFLYDVVA--F** sehhidksyn **ELLF** kwddsqqk **HRLMLFVNNDNPTQAKAELSI** yledyl--sytqa
rlklls **FYNASASKYNKNIDLV** r -----mnk **WNYGVFFVYDVIN--I** ddhylvkkds **ELVF** kwddineee **HQLMLHVNVNEAETVAKEELKL** ienyv--actqp

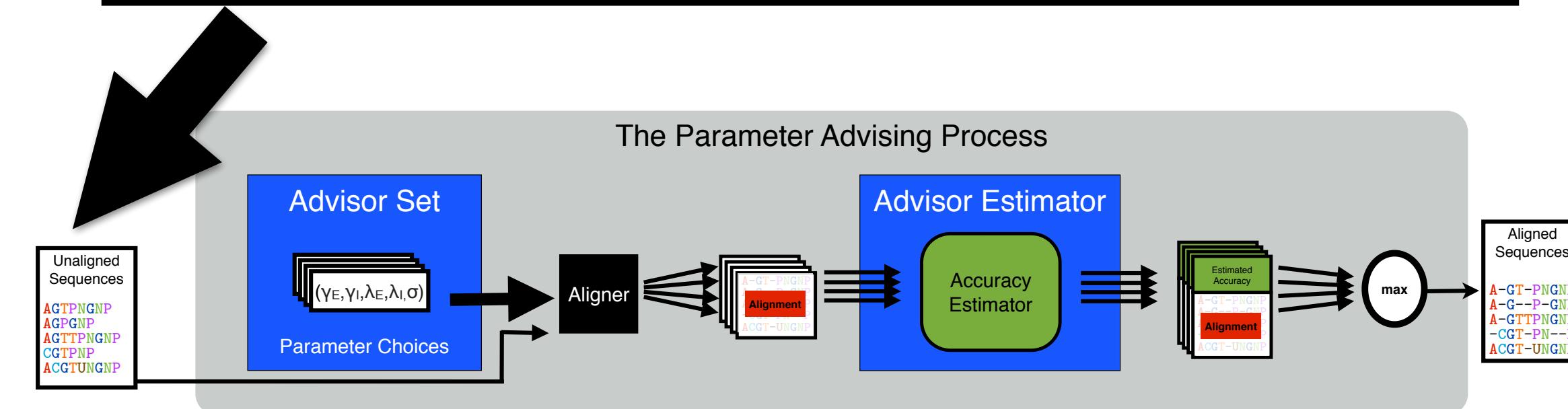
q **MKFGLFFLFDTLAV** yenhfsnngvvldqmsegrfafhkiindafttgcypnnd
- **MKFGNFLLFDTVWL** lehhftefgllldqmskgrfrfydlmkegfnegyiaadne
-----mtk **WNYGVFFFLYDVVA--F** sehhidksyn
-----mnk **WNYGVFFVYDVIN--I** ddhylvkkds



Adaptive local realignment

rkeyag **LYHEVAQAHGVDSQV** q **MKFGLFFLFDTLAV** yenhfsnngvvldqmsegrfafhkiindafttgcypnnd **ELVF** rwddsnaq **HKLTLLVNQNVGEAARAEARV** yleefvresysnt
kkaqlid **LYNEVATEHGYDVTKI** d - **MKFGNFLLFDTVWL** lehhftefgllldqmskgrfrfydlmkegfnegyiaadne **PMIL** swiinthe **HCLSYITSVDHDSNRAKDICRN** flghwy-dsyvna
rielln **HYQAAAAKFNVDIANV** r -----mtk **WNYGVFFFLYDVVA--F** sehhidksyn **ELLF** kwddsqqk **HRLMLFVNNDNPTQAKAELSI** yledyl--sytqa
rlklls **FYNASASKYNKNIDLV** r -----mnk **WNYGVFFVYDVIN--I** ddhylvkkds **ELVF** kwddineee **HQLMLHVNVNEAETVAKEELKL** ienyv--actqp

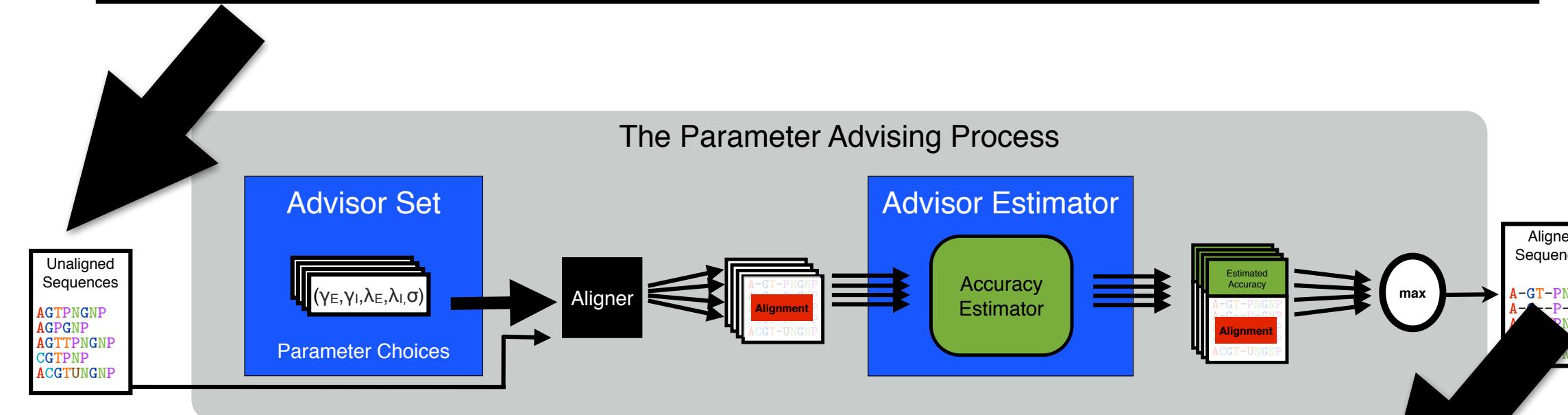
q **MKFGLFFLFDTLAV** yenhfsnngvvldqmsegrfafhkiindafttgcypnnd
- **MKFGNFLLFDTVWL** lehhftefgllldqmskgrfrfydlmkegfnegyiaadne
-----mtk **WNYGVFFFLYDVVA--F** sehhidksyn
-----mnk **WNYGVFFVYDVIN--I** ddhylvkkds



Adaptive local realignment

rkeyag**LYHEVAQAHGVDSQV**rq**MKFGLFFLFDTLAV**yenhsnngvvldqmsegrfafhkiindafattgychpnnd**ELVF**rwddsnaq**HKL****TLLVNQNVDGEAARAEARV**yleefvresysnt
kkaqlid**LYNEVATEHGYDVTKI**d**-MKFGNFLLFDTVWL**lehhftefgl11dqmskgrfrfydlmkegfnegyiaadne**PMIL**swiinthe**HCLSYITSVDHDSNRAKDICRN**flghwy-dsyvna
rielln**HYQAAAAKFNVDIANV**r-----mtk**WNYGVFFFLYDVVA--F**sehhidksyn**ELLF**kwddsqqk**HRLMLFVNNDNPTQAKAELSI**yledyl--sytqa
rlklls**FYNASASKYNKNIDLV**r-----mnk**WNYGVFFVYDVIN--I**ddhylvkkds**ELVF**kwddineee**HQLMLHVNVNEAETVAKEELKL**yienyv--actqp

q**MKFGLFFLFDTLAV**yenhsnngvvldqmsegrfafhkiindafattgychpnnd
-**MKFGNFLLFDTVWL**lehhftefgl11dqmskgrfrfydlmkegfnegyiaadne
-----mtk**WNYGVFFFLYDVVA--F**sehhidksyn
-----mnk**WNYGVFFVYDVIN--I**ddhylvkkds

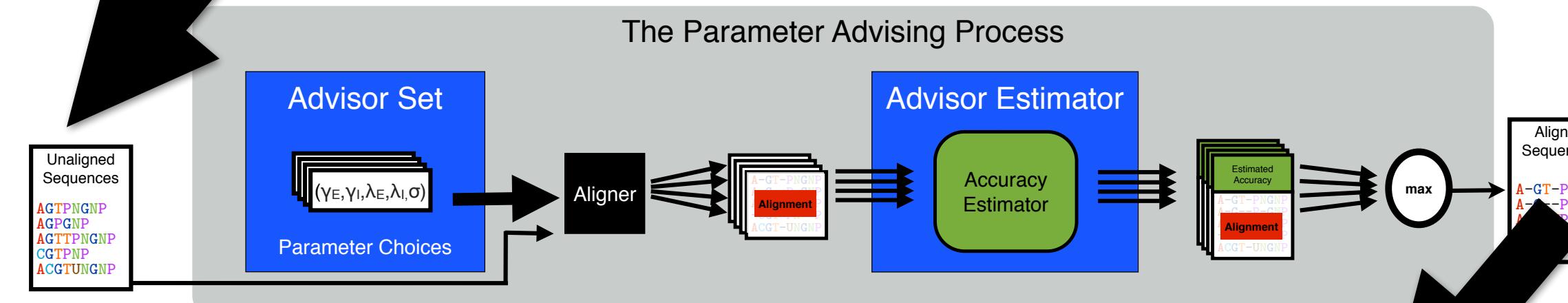


---q**MKFGLFFLFDTLAV**yenhsnngvvldqmsegrfafhkiindafattgychpnnd
---**MKFGNFLLFDTVWL**lehhftefgl11dqmskgrfrfydlmkegfnegyiaadne
mtk**WNYGVFFFLYDVVAF**sehhidksyn-----
mnk**WNYGVFFVYDVINI**ddhylvkkds-----

Adaptive local realignment

rkeyag **LYHEVAQAHGVDSQV** r q **MKFGLFFLFDTLAV** yenhfsnngvvldqmsegrfafhkiindafttgychpnnd **ELVF** rwddsnaq **HKLTLLVNQNVGEAARAEARV** yleefvresysnt
kkaql d **LYNEVATEHGYDVTKI** d - **MKFGNFLLFDTVWL** lehhftefgllldqmskgrfrfydlmkegfnegyiaadne **EMIL** swiinthe **HCLSYITSVDHDSNRAKDICRN** flghwy-dsyvna
rielln **HYQAAAAKFNVDIANV** r mtk **WNYGVFFFLYDVVA--F** sehhidksyn **ELLF** kwddsqqk **HRLMLFVNNDNPTQAKAELSI** yledyl--sytqa
rlklls **FYNASASKYNKNIDLV** r mnk **WNYGVFFVYDVIN--I** ddhylvkkds **ELVF** kwddineee **HQLMLHVNVNEAETVAKEELKL** ienyv--actqp

q **MKFGLFFLFDTLAV** yenhfsnngvvldqmsegrfafhkiindafttgychpnnd
- **MKFGNFLLFDTVWL** lehhftefgllldqmskgrfrfydlmkegfnegyiaadne
-----mtk **WNYGVFFFLYDVVA--F** sehhidksyn
-----mnk **WNYGVFFVYDVIN--I** ddhylvkkds



---q **MKFGLFFLFDTLAV** yenhfsnngvvldqmsegrfafhkiindafttgychpnnd
--- **MKFGNFLLFDTVWL** lehhftefgllldqmskgrfrfydlmkegfnegyiaadne
mtk **WNYGVFFFLYDVVAF** sehhidksyn-----
mnk **WNYGVFFVYDVINI** ddhylvkkds-----

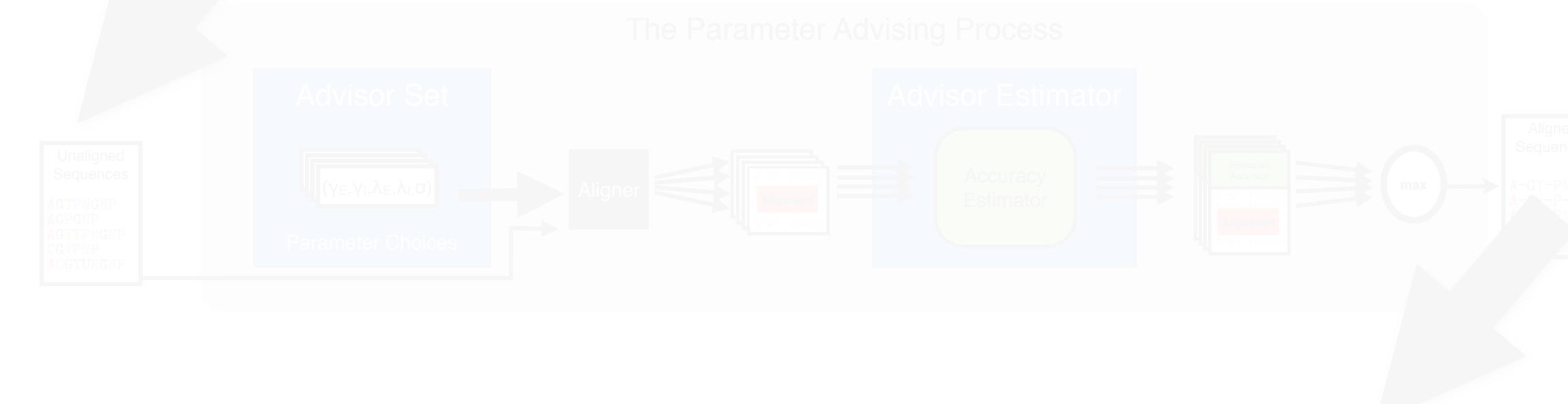
rkeyag **LYHEVAQAHGVDSQV** r --q **MKFGLFFLFDTLAV** yenhfsnngvvldqmsegrfafhkiindafttgychpnnd **ELVF** rwddsnaq **HKLTLLVNQNVGEAARAEARV** yleefvresysnt
kkaql d **LYNEVATEHGYDVTKI** d -- **MKFGNFLLFDTVWL** lehhftefgllldqmskgrfrfydlmkegfnegyiaadne **EMIL** swiinthe **HCLSYITSVDHDSNRAKDICRN** flghwy-dsyvna
rielln **HYQAAAAKFNVDIANV** r mtk **WNYGVFFFLYDVVAF** sehhidksyn----- **ELLF** kwddsqqk **HRLMLFVNNDNPTQAKAELSI** yledyl--sytqa
rlklls **FYNASASKYNKNIDLV** r mnk **WNYGVFFVYDVINI** ddhylvkkds----- **ELVF** kwddineee **HQLMLHVNVNEAETVAKEELKL** ienyv--actqp

Adaptive local realignment

rkeyagLYHEVAQAHGVDSQVr qMKFGLFFLFDTLAVyenhsnnngvvld
kkaql dLYNEVATEHGYDVTKId -MKFGNFLLFDTVWLlehhftefglllde
riellnHYQAAAAKFNDIANVr mtk wNYGVFFFLYDVVAAF sehhidksyn
rlkllsFYNASASKYNKNIDLVr mnk wNYGVFFVYDVINI ddhylvkkds

Accuracy
33%

q**MKFGLFFLFDTLAV**yenhsnnngvvldqmsegrfafhkiindafttgychpnnd
-**MKFGNFLLFDTVWL**lehhftefgllldqmskgrfrfydlmkegfnegyiaadne
-----mtk **WNYGVFFFLYDVVA--F** sehhidksyn
-----mnk **WNYGVFFVYDVIN--I** ddhylvkkds



---q**MKFGLFFLFDTLAV**yenhsnnngvvldqmsegrfafhkiindafttgychpnnd
---**MKFGNFLLFDTVWL**lehhftefgllldqmskgrfrfydlmkegfnegyiaadne
mtk **WNYGVFFFLYDVVAAF** sehhidksyn-----
mnk **WNYGVFFVYDVINI** ddhylvkkds-----

Accuracy
100%

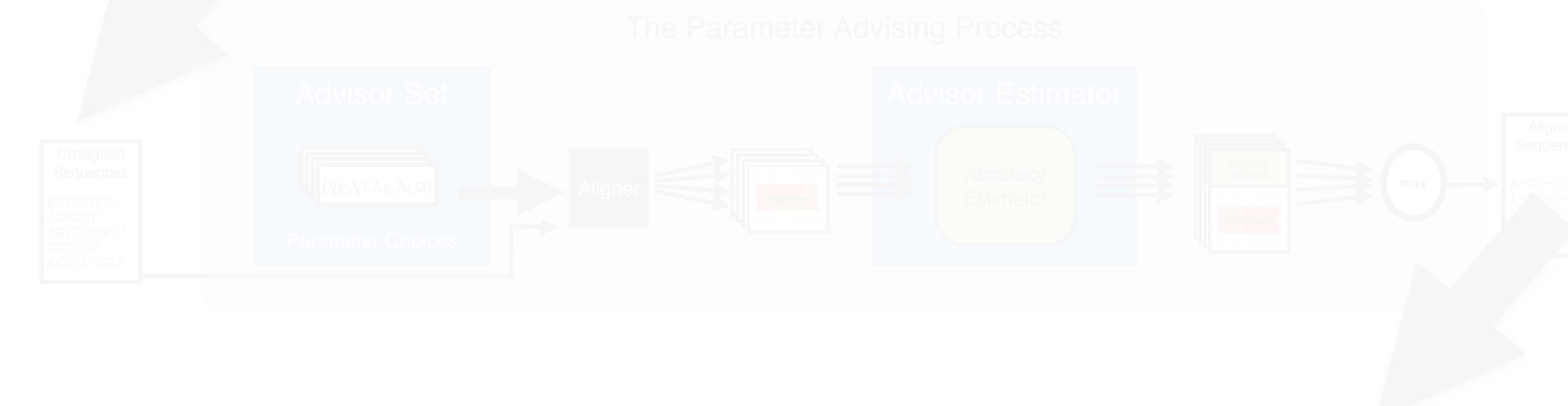
rkeyagLYHEVAQAHGVDSQVr --qMKFGLFFLFDTLAVyenhsnnngvvld
kkaql dLYNEVATEHGYDVTKId --MKFGNFLLFDTVWLlehhftefglllde
riellnHYQAAAAKFNDIANVr mtk wNYGVFFFLYDVVAAF sehhidksyn--
rlkllsFYNASASKYNKNIDLVr mnk wNYGVFFVYDVINI ddhylvkkds--

ndPLVFrwddsnaq**HKLTLVNVNQVDGEAARAEARV**yleefvresysnt
dneBMLswiinthe**HCLSYITSVDHDSNRAKDICRN**flghwy-dsyvna
ynELLFkwddsqqk**HRLMLFVNNDNPTQAKAELSI**yledyl--sytqa
dsPLVfkwdidinee**HQLMLHVNNEAETVAKEELKL**yienyv--actqp

Adaptive local realignment

rkeyag **LYHEVAQAHGVDSQV** r **MKFGLFFLFDTLAV** yenhfsnnvvldqmsegrfafhkiindafttgcypnnd **ELVF** rwddsnaq **HKLTLLVNQNVGEAARAEARV** yleefvresysnt
kkaql d **LYNEVATEHGYDVTKI** d **-MKFGNFLLFDTVWL** lehhftefgllldqmskgrfrfydlmkegfnegyiaadne **EMIL** swiinthe **HCLSYITSVDHDSNRACKDICRN** flghwy-dsyvna
rielln **HYQAAAAKFNVDIANV** r mtk **WNYGVFFFLYDVVA--F** sehhidksyn **ELLF** kwddsqqk **HRLMLFVNNDNPTQAKAELSI** yledyl--sytqa
rlklls **FYNASASKYNKNIDLV** r mnk **WNYGVFFVYDVIN--I** ddhylvkkds **ELVF** kwddineee **HQLMLHVNVNEAETVAKEELKL** ienyv--actqp

Accuracy
80%



Accuracy
100%

rkeyag **LYHEVAQAHGVDSQV** r **-qMKFGLFFLFDTLAV** yenhfsnnvvldqmsegrfafhkiindafttgcypnnd **ELVF** rwddsnaq **HKLTLLVNQNVGEAARAEARV** yleefvresysnt
kkaql d **LYNEVATEHGYDVTKI** d **--MKFGNFLLFDTVWL** lehhftefgllldqmskgrfrfydlmkegfnegyiaadne **EMIL** swiinthe **HCLSYITSVDHDSNRACKDICRN** flghwy-dsyvna
rielln **HYQAAAAKFNVDIANV** r mtk **WNYGVFFFLYDVVAF** sehhidksyn ----- **ELLF** kwddsqqk **HRLMLFVNNDNPTQAKAELSI** yledyl--sytqa
rlklls **FYNASASKYNKNIDLV** r mnk **WNYGVFFVYDVINI** ddhylvkkds ----- **ELVF** kwddineee **HQLMLHVNVNEAETVAKEELKL** ienyv--actqp

Adaptive local realignment

Adaptive local realignment takes as **input**

Adaptive local realignment

Adaptive local realignment takes as **input**

- an **initial alignment**,

Adaptive local realignment

Adaptive local realignment takes as **input**

- an **initial alignment**,
- an **advisor set**,

Adaptive local realignment

Adaptive local realignment takes as **input**

- an **initial alignment**,
- an **advisor set**,
- **window sizes**, and

Adaptive local realignment

Adaptive local realignment takes as **input**

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

Adaptive local realignment

Adaptive local realignment takes as **input**

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

And **outputs**

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**

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

Experimental results

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

- with the Opal aligner and Facet estimator,

Experimental results

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

- with the Opal **aligner** and Facet **estimator**,
- on over 800 **benchmarks** from BENCH and PALI,

Experimental results

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

- with the Opal **and Facet **,****
- on over 800 **from BENCH and PALI,**
- using 12-fold **.**

Experimental results

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

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.

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**.

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.

Experimental results

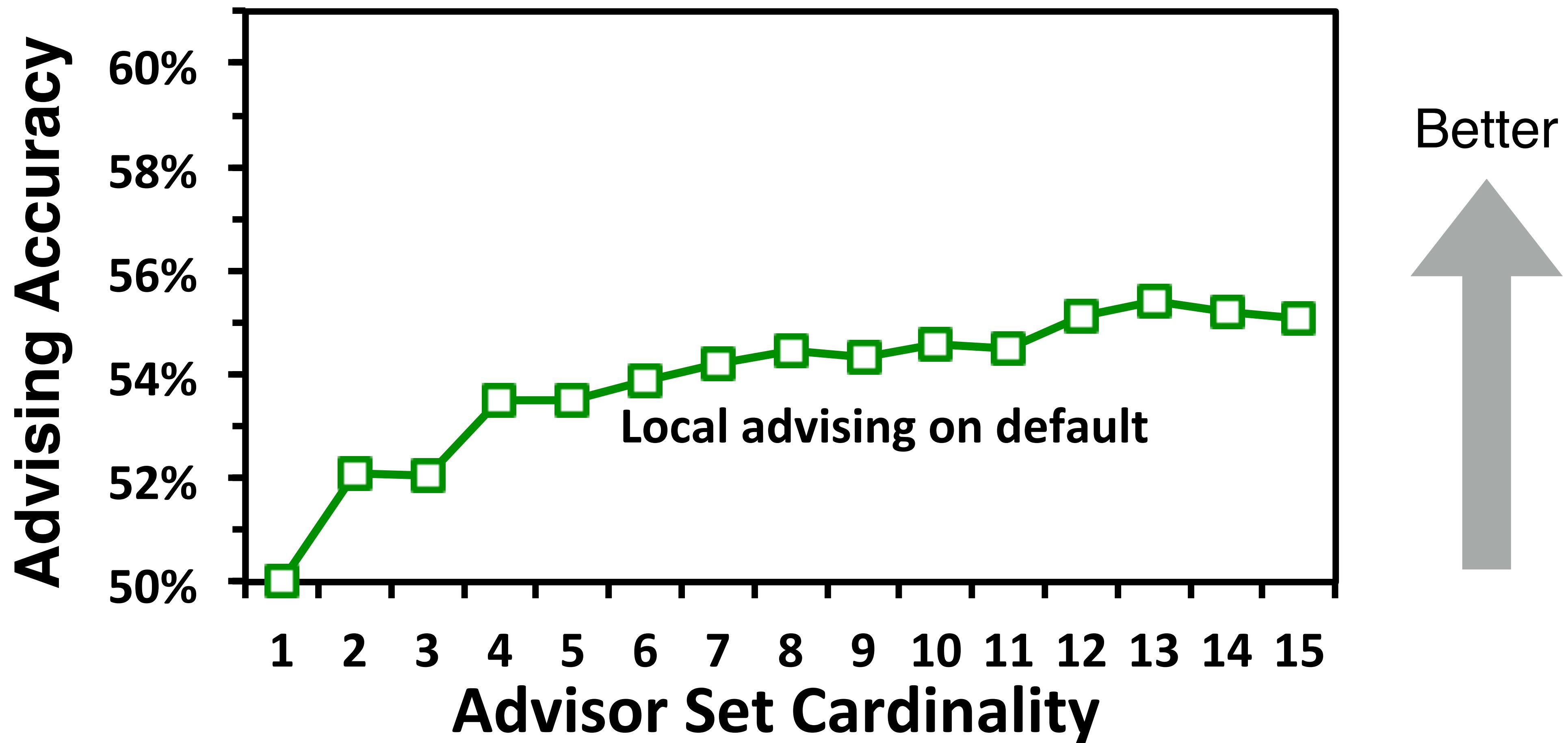
We correct for the **bias** in over-representation of easy-to-align benchmarks.

- 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%.

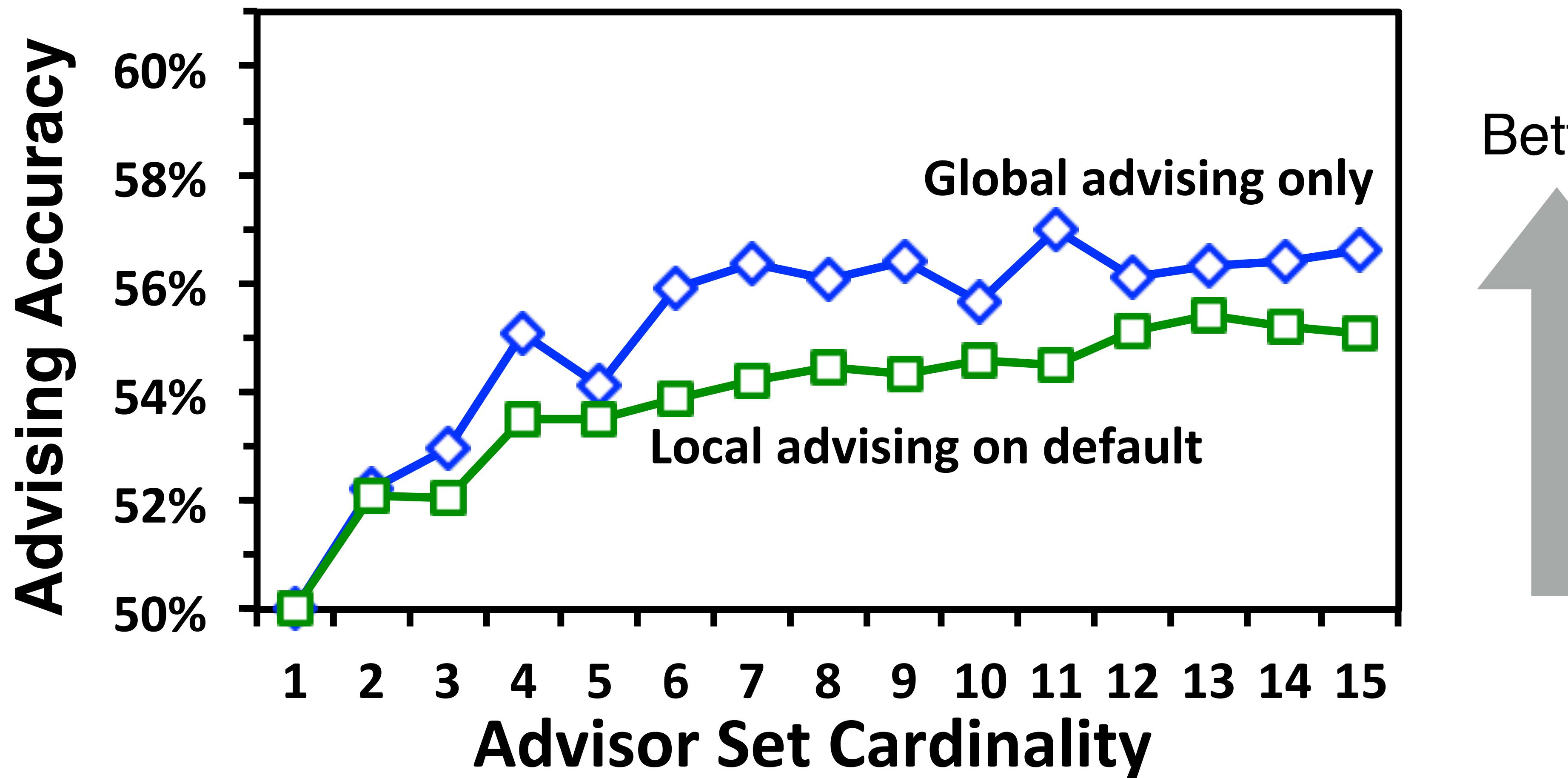
Experimental results

Average accuracy of advisors versus set cardinality



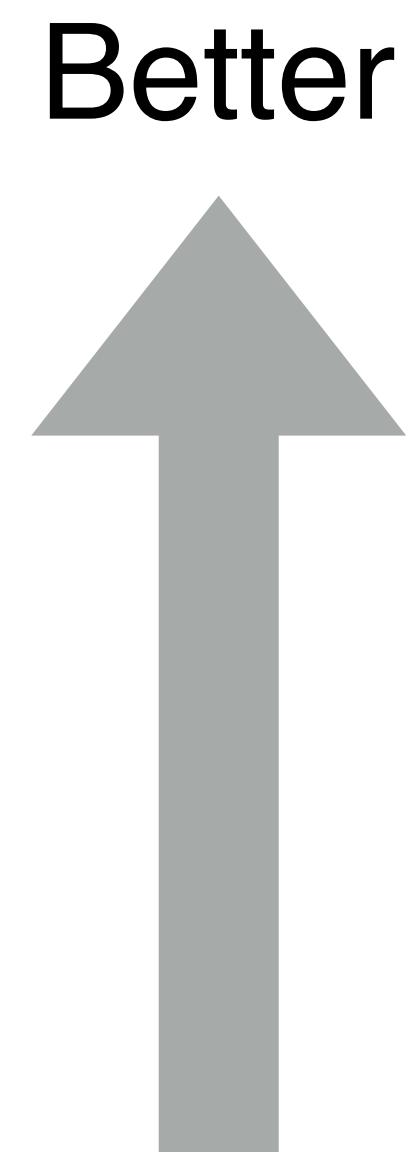
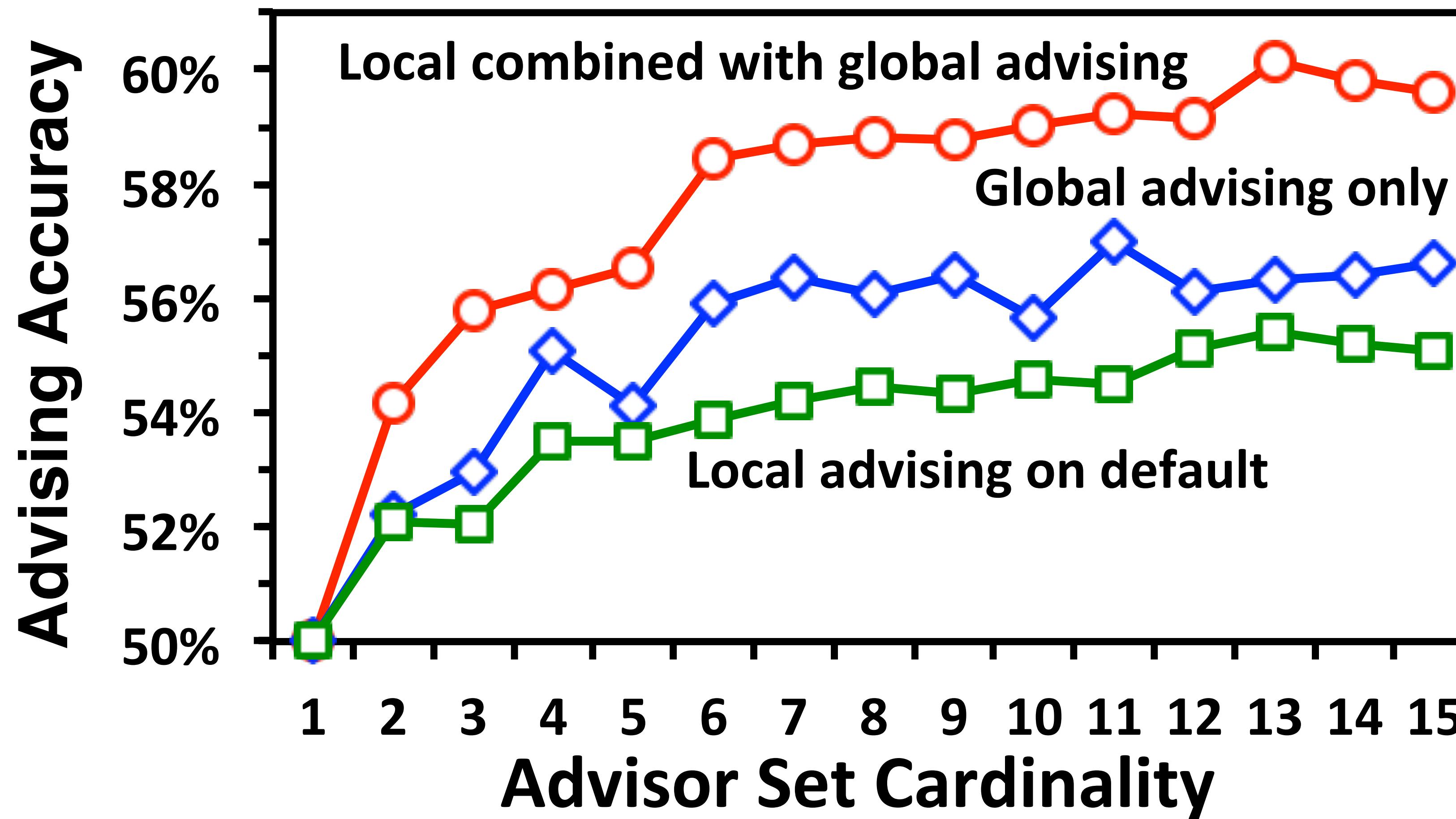
Experimental results

Average accuracy of advisors versus set cardinality



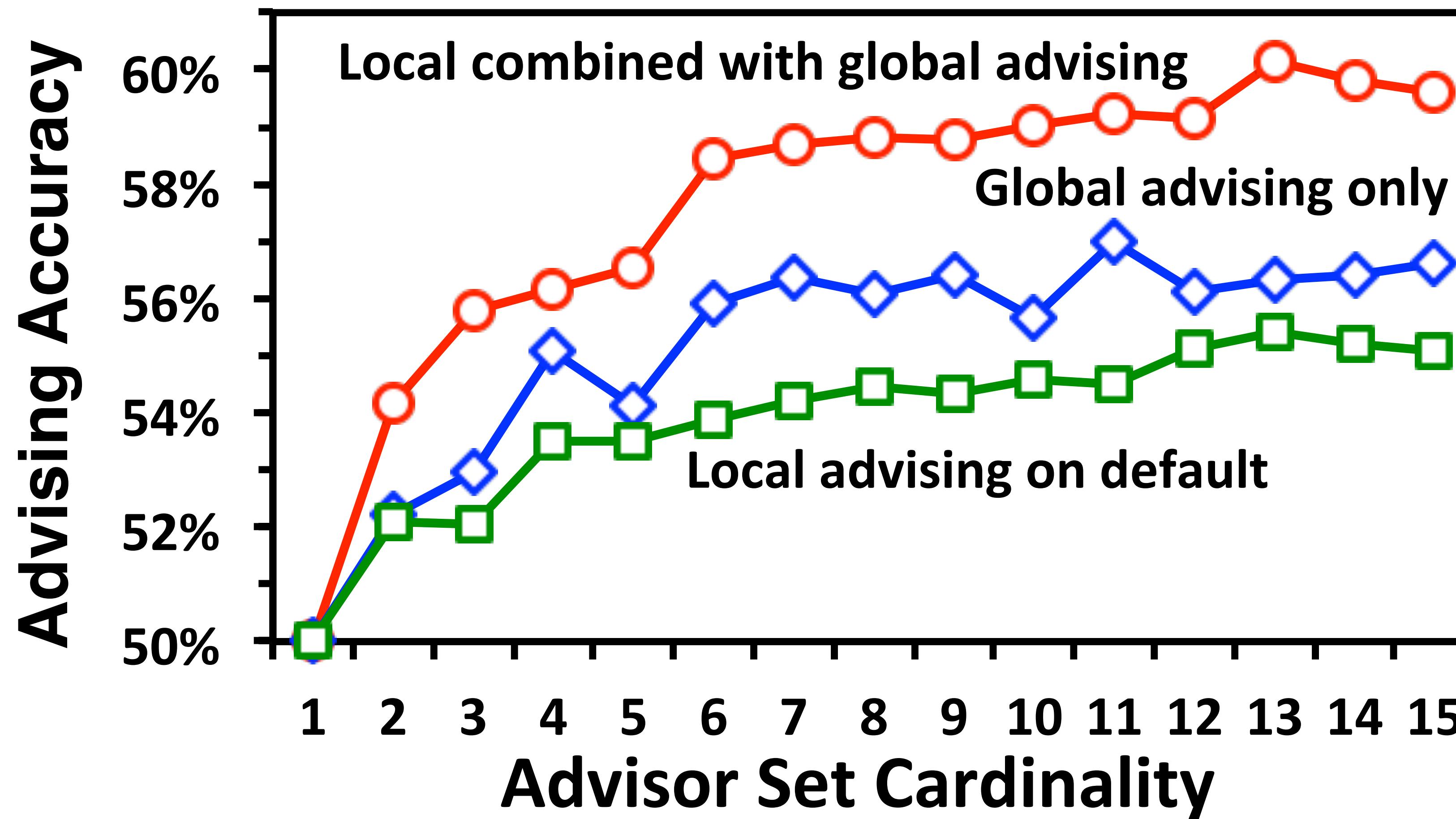
Experimental results

Average accuracy of advisors versus set cardinality



Experimental results

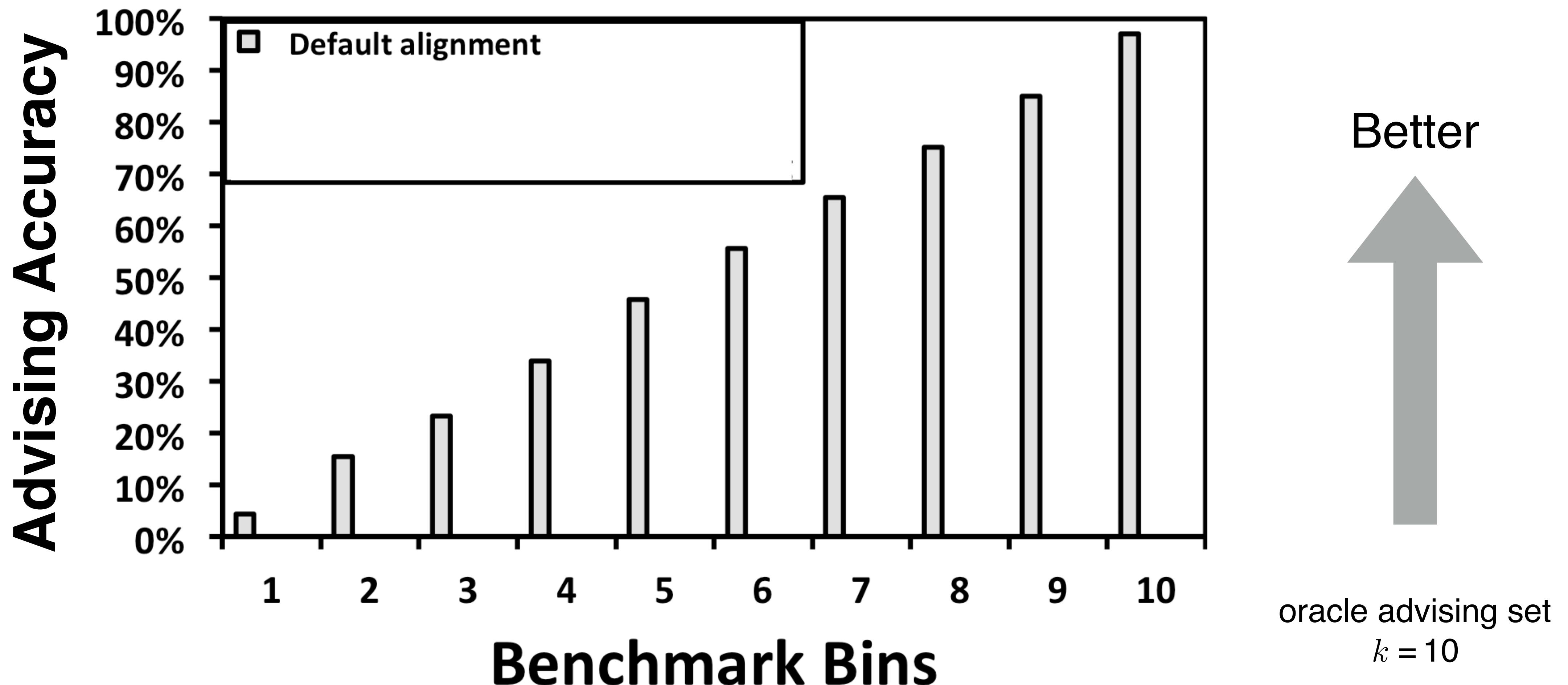
Average accuracy of advisors versus set cardinality



Boosts accuracy by 9% over default parameter choice

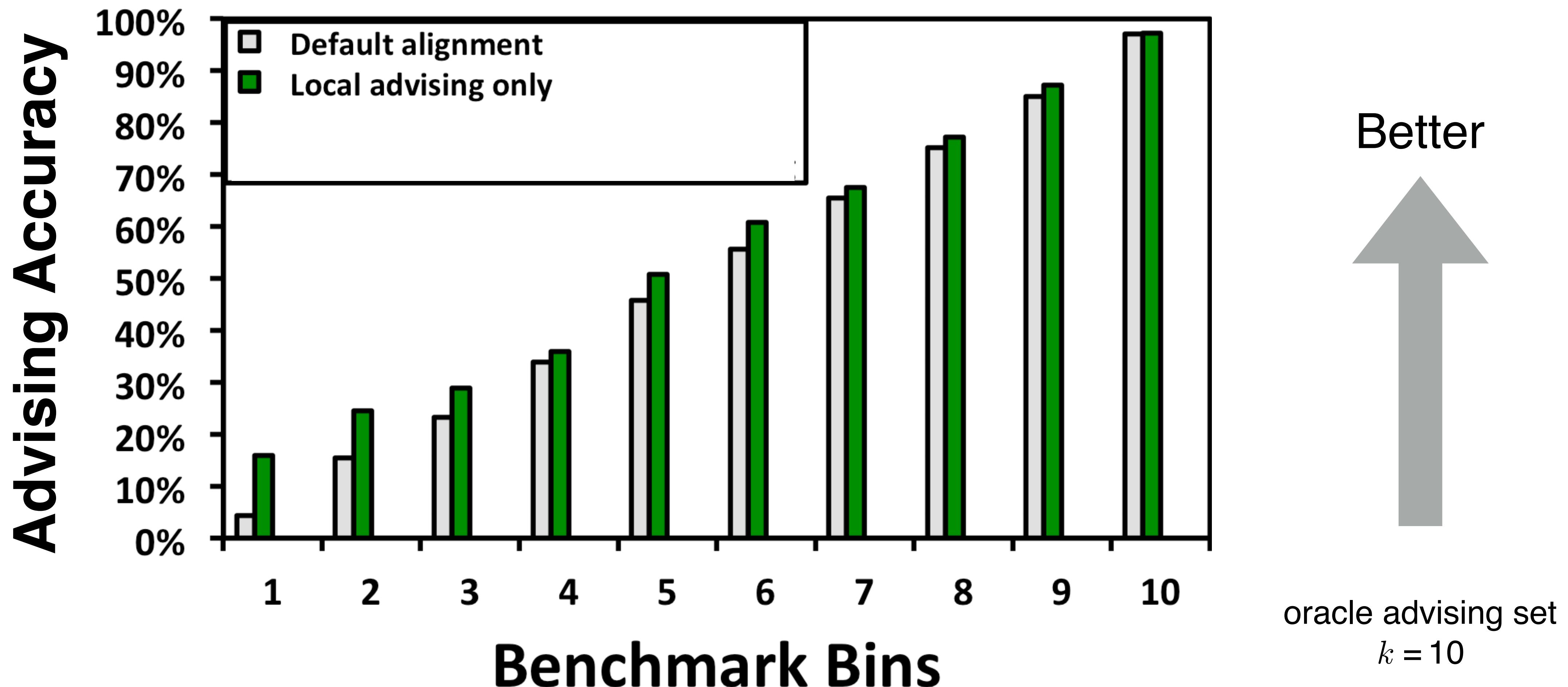
Experimental results

Average accuracy of advisors within difficulty bins



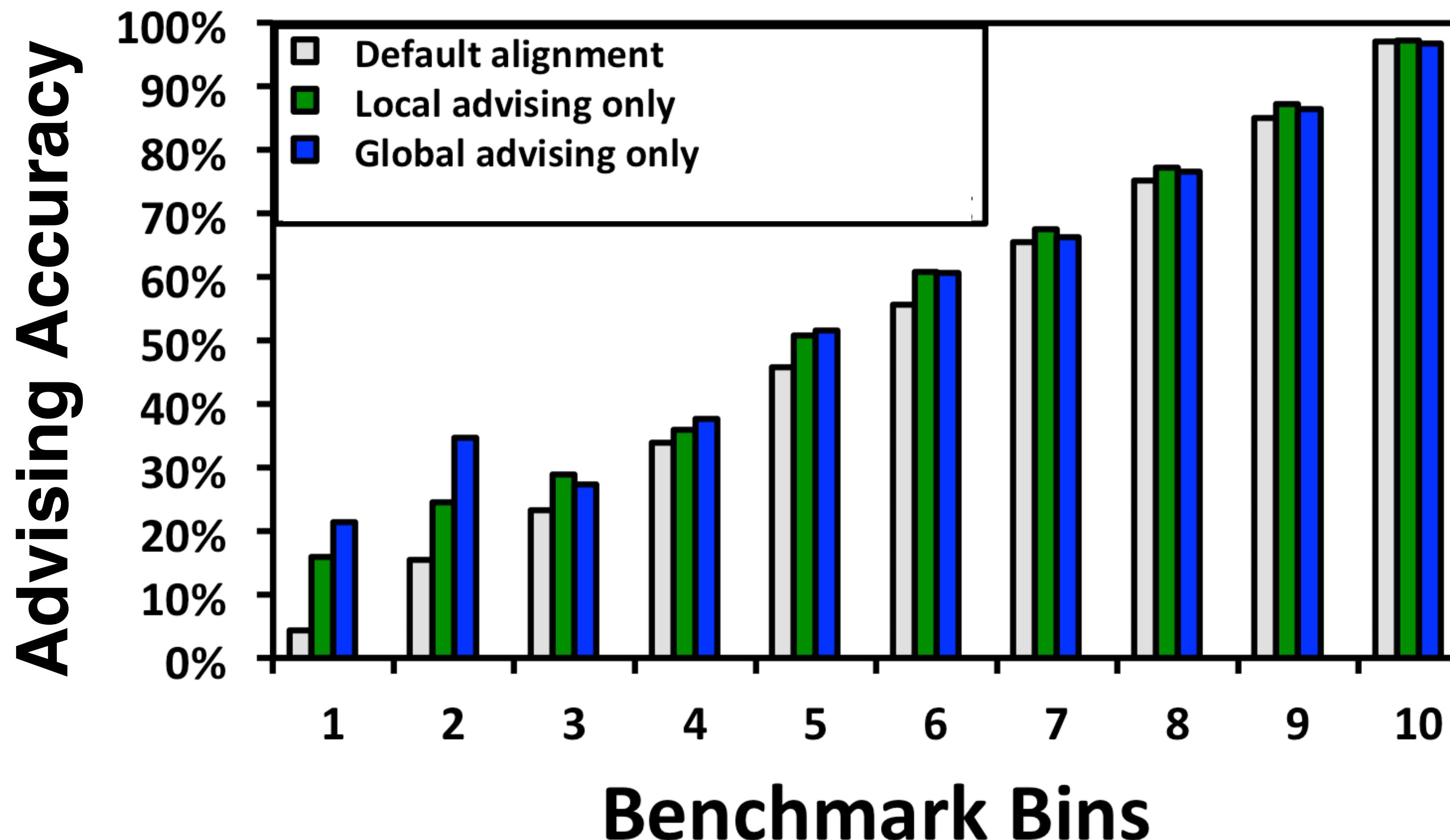
Experimental results

Average accuracy of advisors within difficulty bins

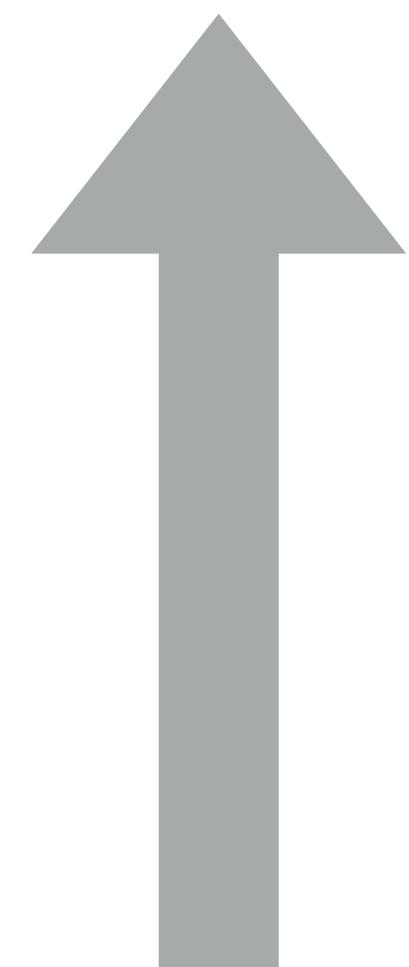


Experimental results

Average accuracy of advisors within difficulty bins



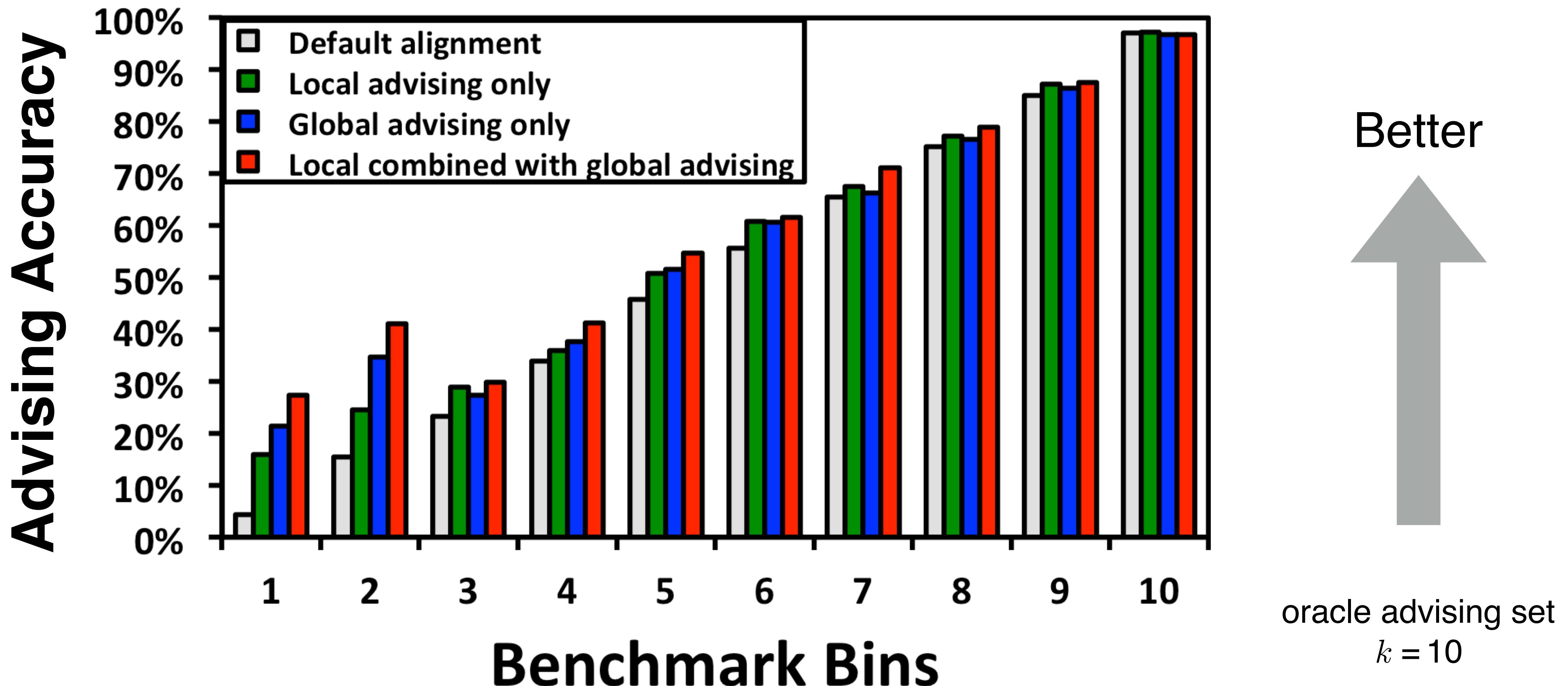
Better



oracle advising set
 $k = 10$

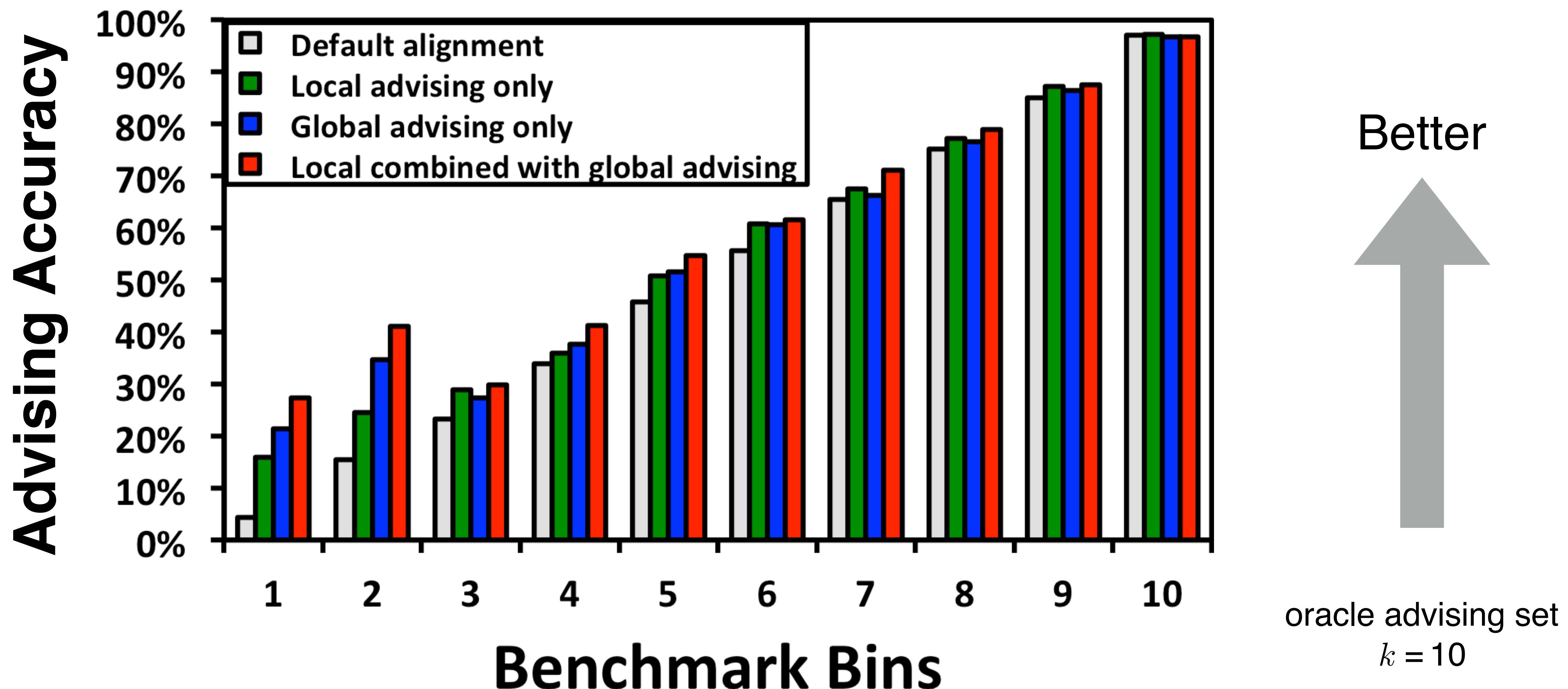
Experimental results

Average accuracy of advisors within difficulty bins



Experimental results

Average accuracy of advisors within difficulty bins



Boosts accuracy on the hardest bins by almost 26%

Conclusions

Local adaptive realignment:

Conclusions

Local adaptive realignment:

- Uses **diverse** parameter choices for heterogeneous proteins.

Conclusions

Local adaptive realignment:

- Uses **diverse** parameter choices for heterogeneous proteins.
- Facet estimator identifies misaligned **regions**.

Conclusions

Local adaptive realignment:

- Uses **diverse** parameter choices for heterogeneous proteins.
- Facet estimator identifies misaligned **regions**.
- Misaligned regions are polished with **parameter advising**.

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:

Further research

Promising directions include:

- Finding advisor sets tuned for local realignment.

Further research

Promising directions include:

- Finding **advisor sets** tuned for **local realignment**.
- Improving the **estimator** by training on **local examples**.

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

Carl Kingsford

Funding

University of Arizona

- NSF IGERT in Genomics Grant DGE-0654435
- NSF Grant IIS-1217886

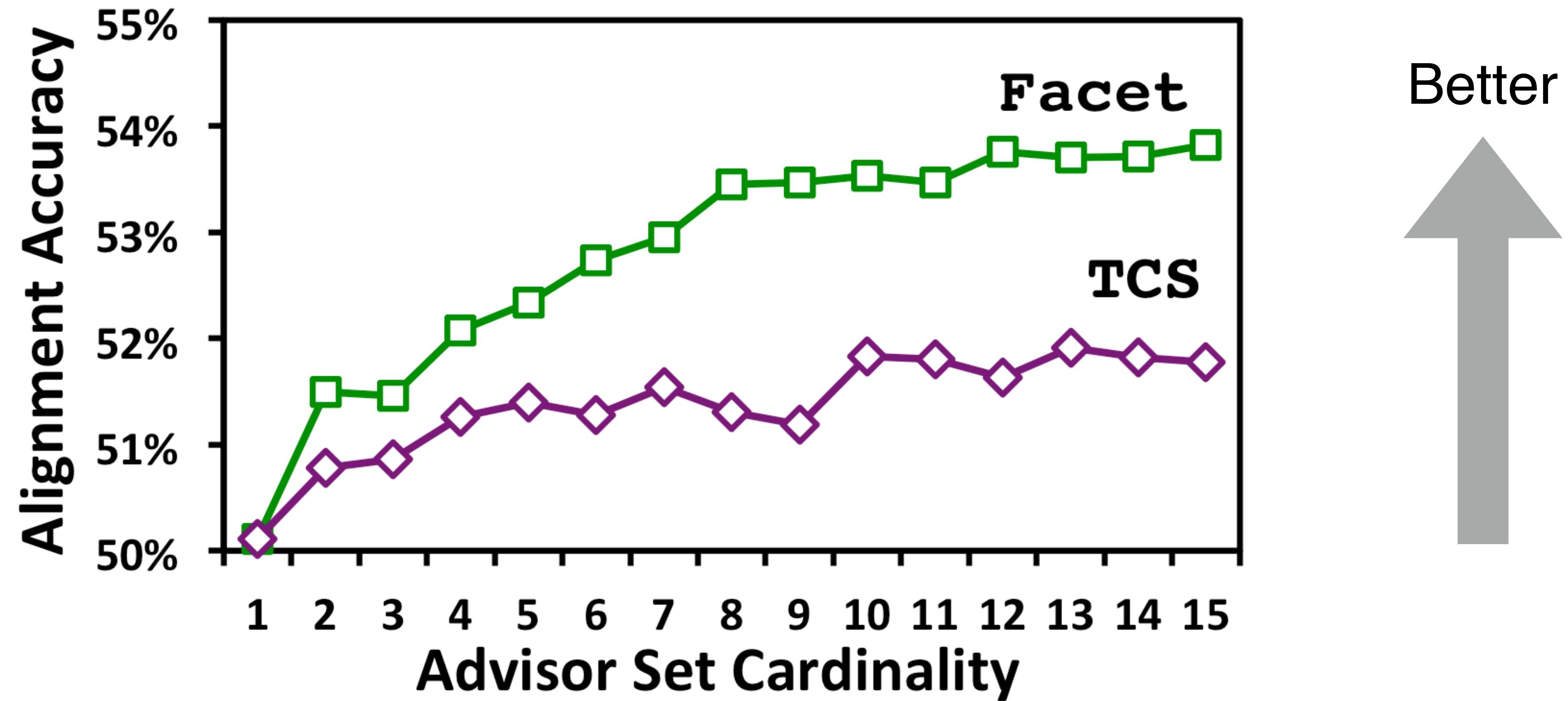
Carnegie Mellon University

- NSF Grants CCF-1256087 & CCF-131999
- NIH Grant R01HG007104
- Gordon and Betty Moore Foundation Grant GBMF4554



Experimental results

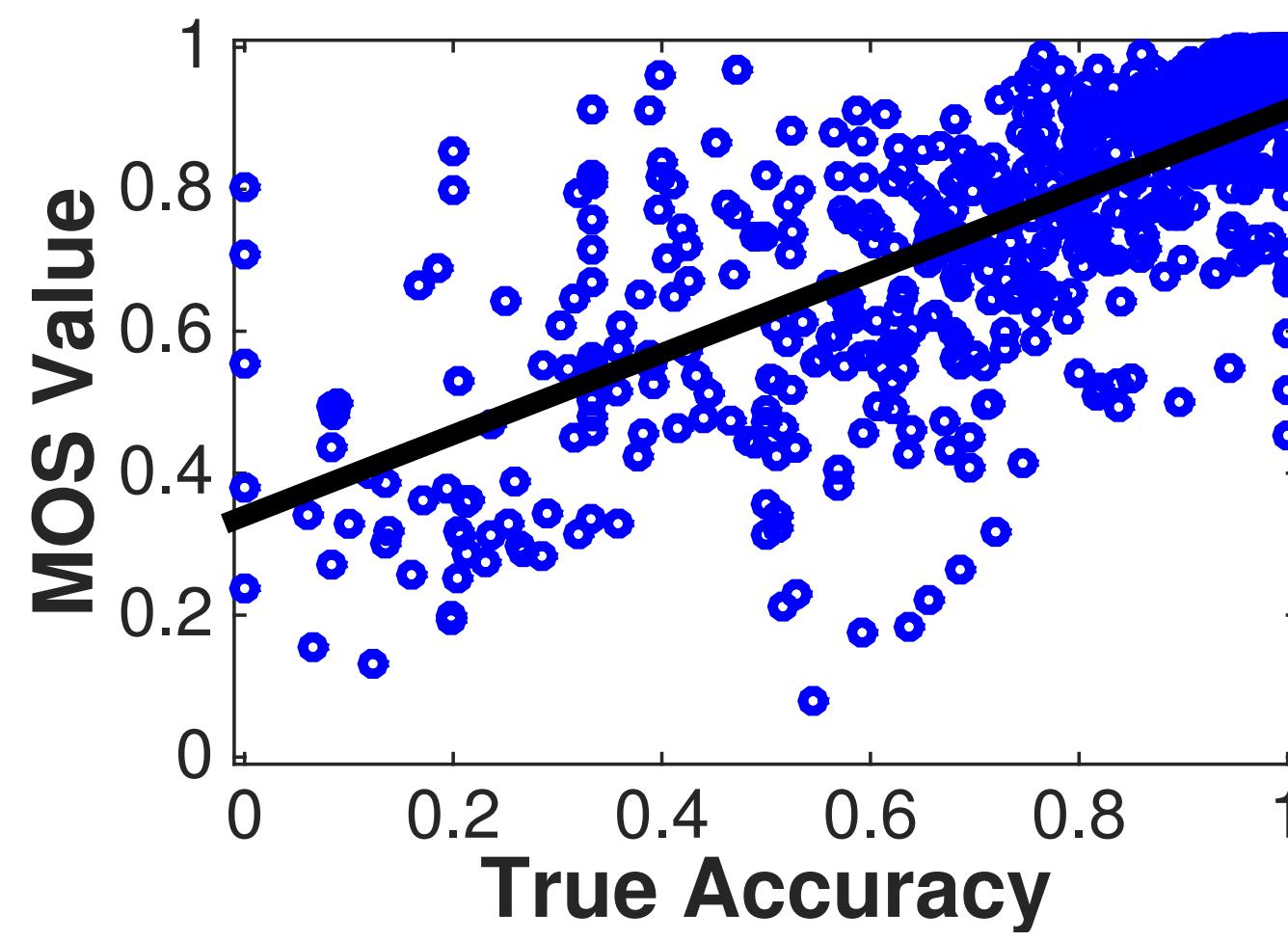
Average accuracy of advisors versus set cardinality



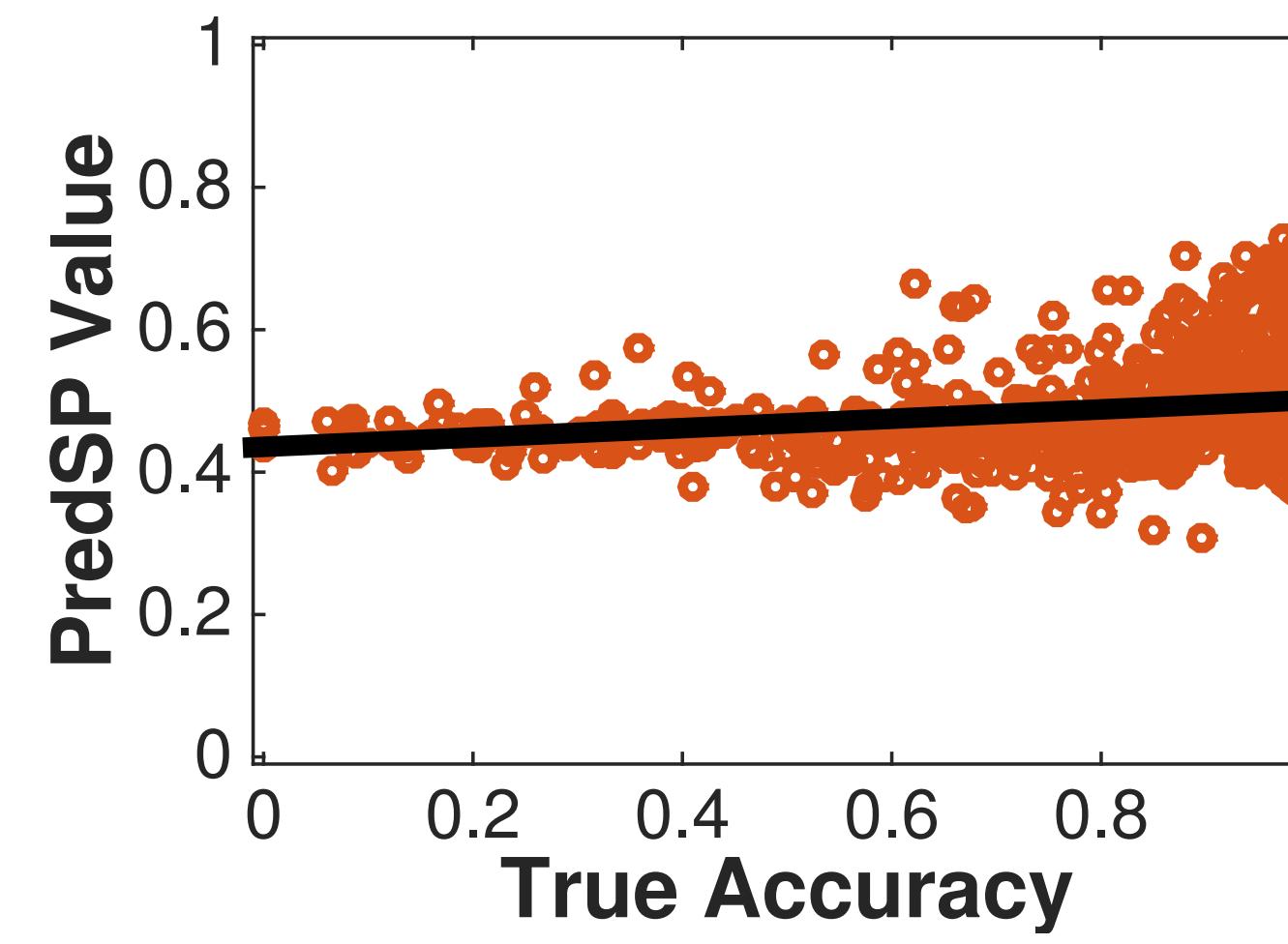
Facet outperforms TCS for adaptive local realignment

Accuracy estimation

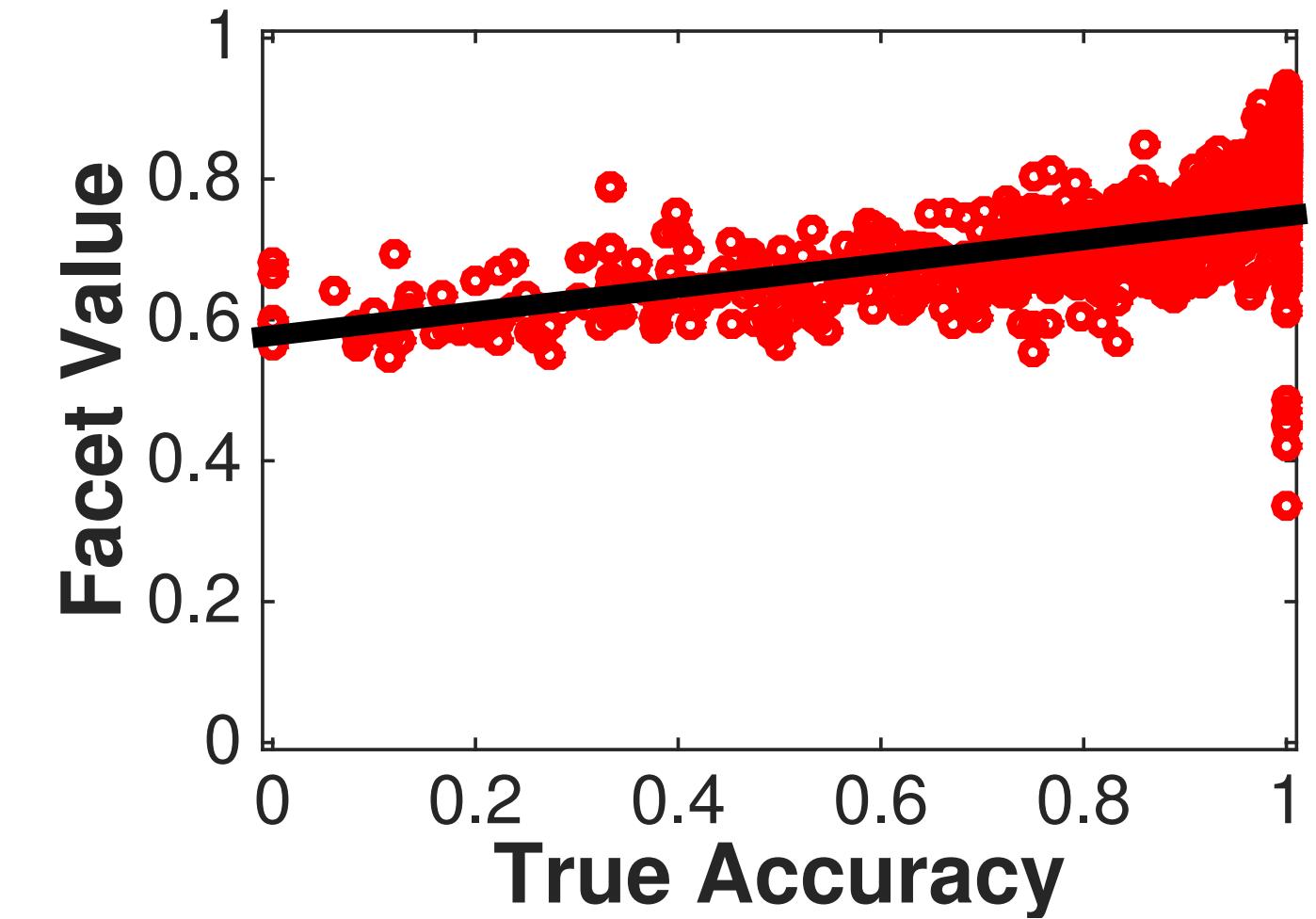
For parameter advising, an estimator should have high **slope** and low **spread**.



high slope,
high spread



low slope,
low spread



medium slope,
low spread

Facet's slope and spread is **best for advising**