

Accurate Identification of Transcription Regulatory Sequences and Genes in Coronaviruses

Chuanyi Zhang^{†,1} Palash Sashittal,^{†‡,2} Michael Xiang,² Yichi Zhang,² Ayesha Kazi,² and Mohammed El-Kebir ^{*,2}

¹Department of Electrical & Computer Engineering, University of Illinois at Urbana-Champaign, Urbana, IL, USA

²Department of Computer Science, University of Illinois at Urbana-Champaign, Urbana, IL, USA

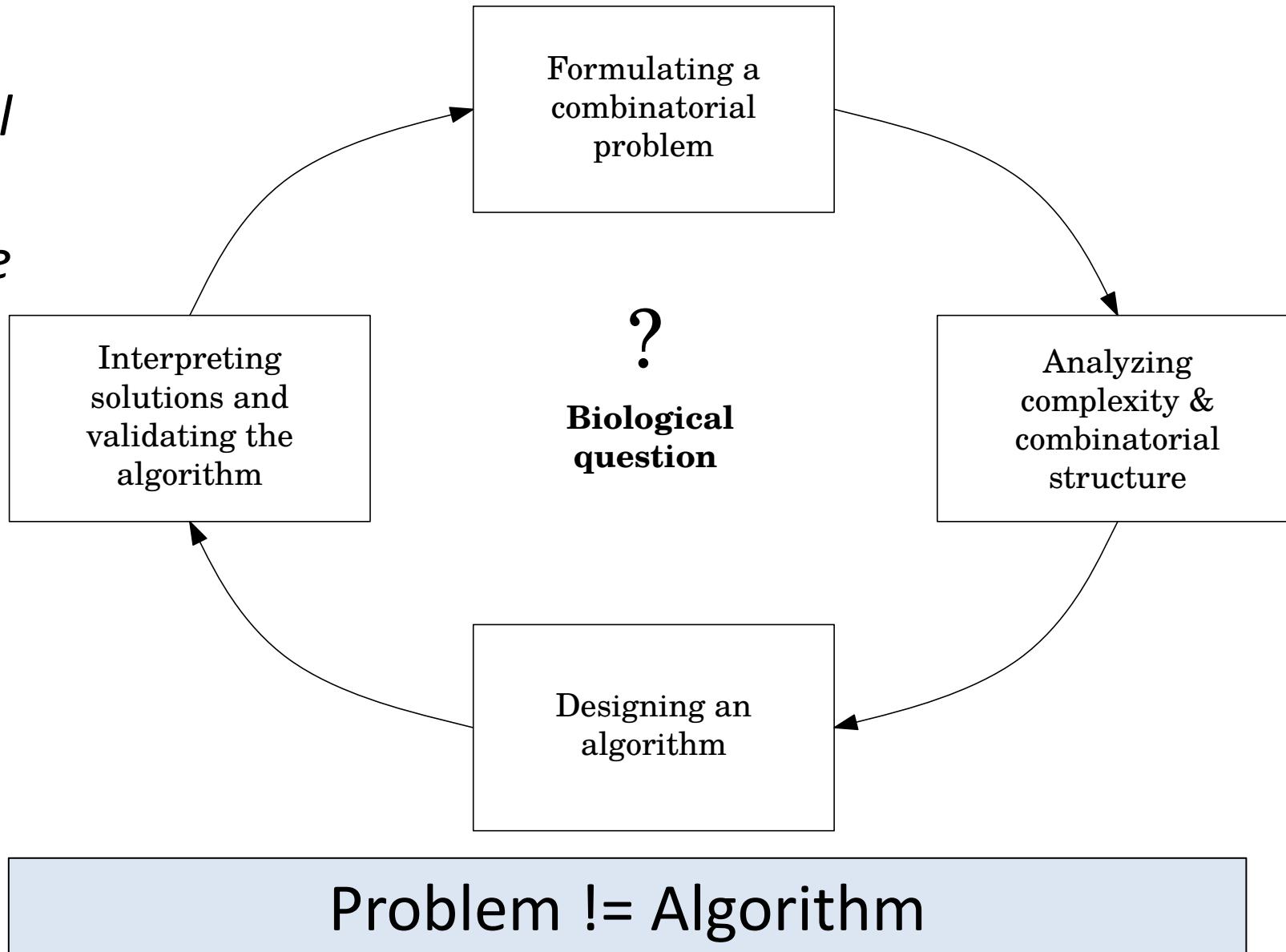
[†]These authors share joint first authorship.

^{*}Present address: Department of Computer Science, Princeton University, Princeton, NJ, USA

Research Statement & Approach

Lab focus:

Application of combinatorial optimization techniques to answers questions and solve problems in biology.



Accurate Identification of Transcription Regulatory Sequences and Genes in Coronaviruses

Chuanyi Zhang^{†,1} Palash Sashittal,^{†‡,2} Michael Xiang,² Yichi Zhang,² Ayesha Kazi,² and Mohammed El-Kebir *,²

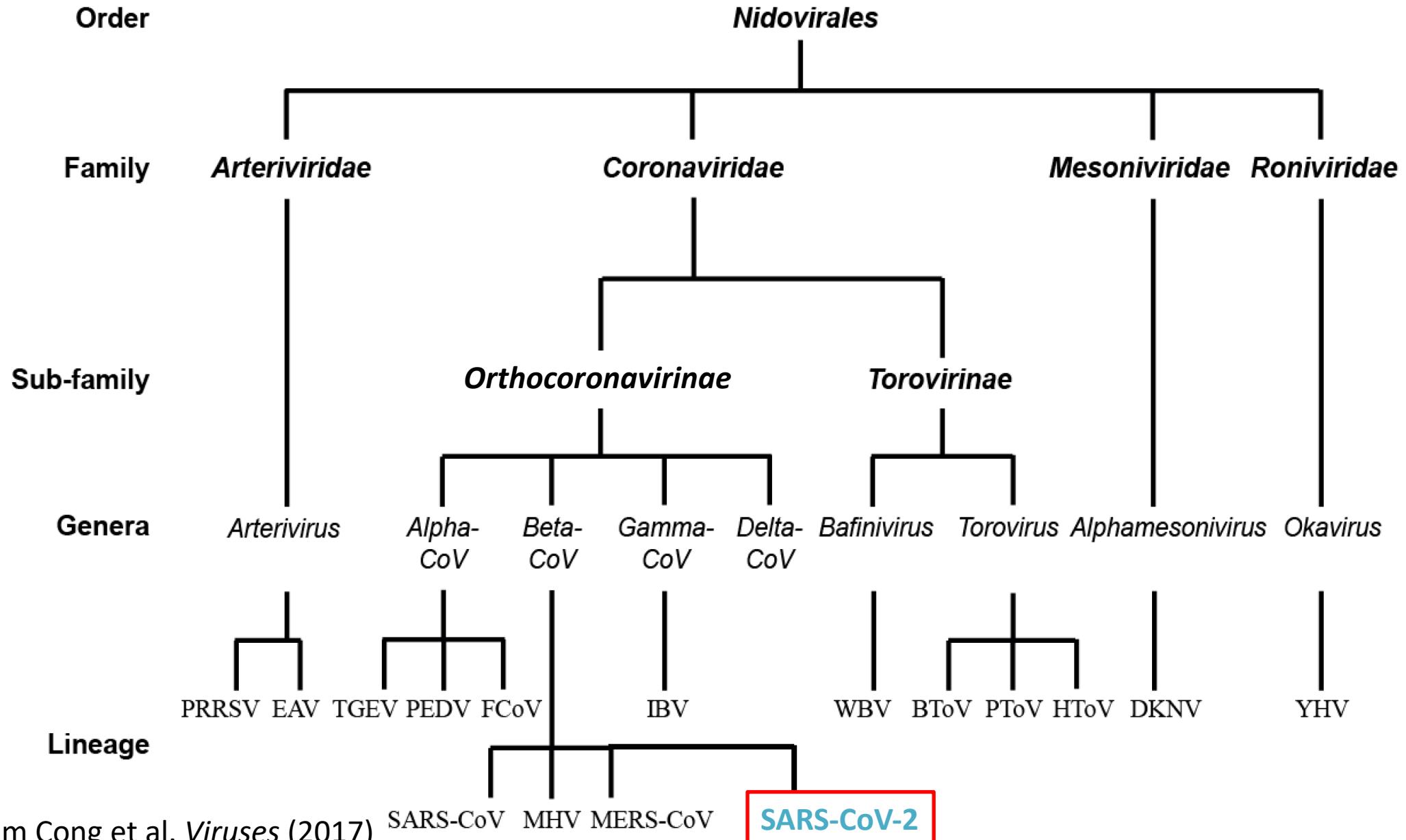
¹Department of Electrical & Computer Engineering, University of Illinois at Urbana-Champaign, Urbana, IL, USA

²Department of Computer Science, University of Illinois at Urbana-Champaign, Urbana, IL, USA

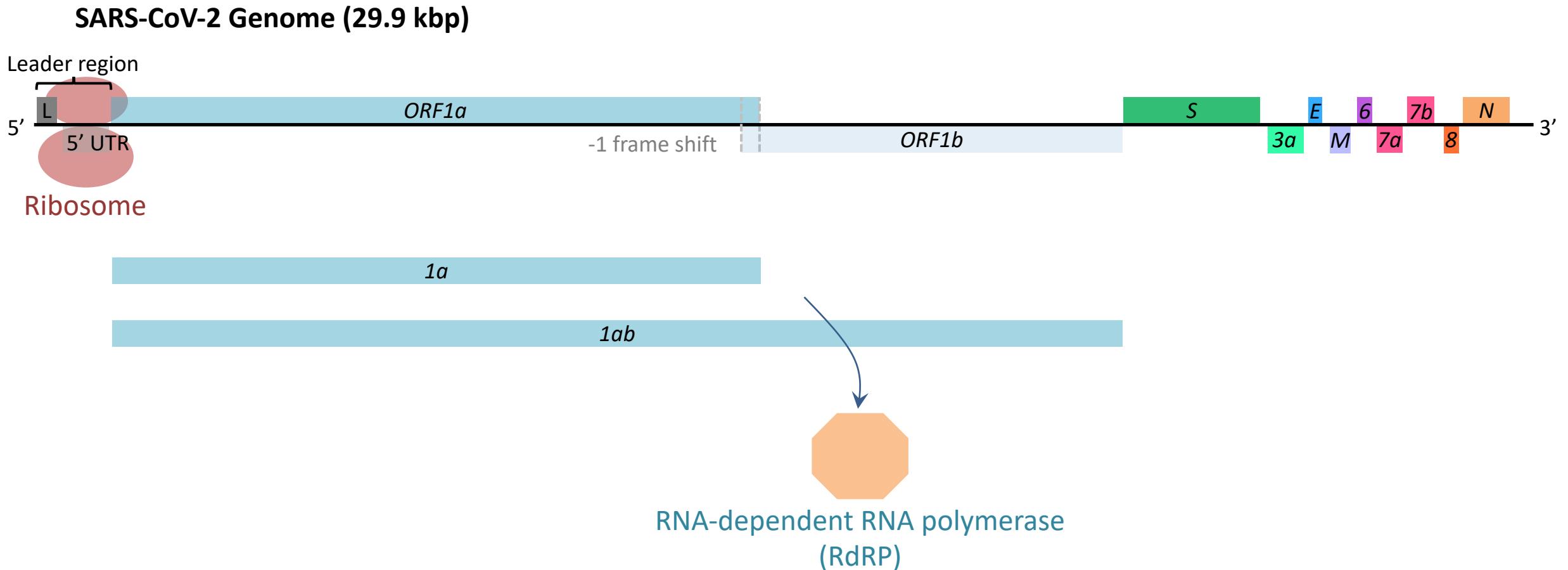
[†]These authors share joint first authorship.

[‡]Present address: Department of Computer Science, Princeton University, Princeton, NJ, USA

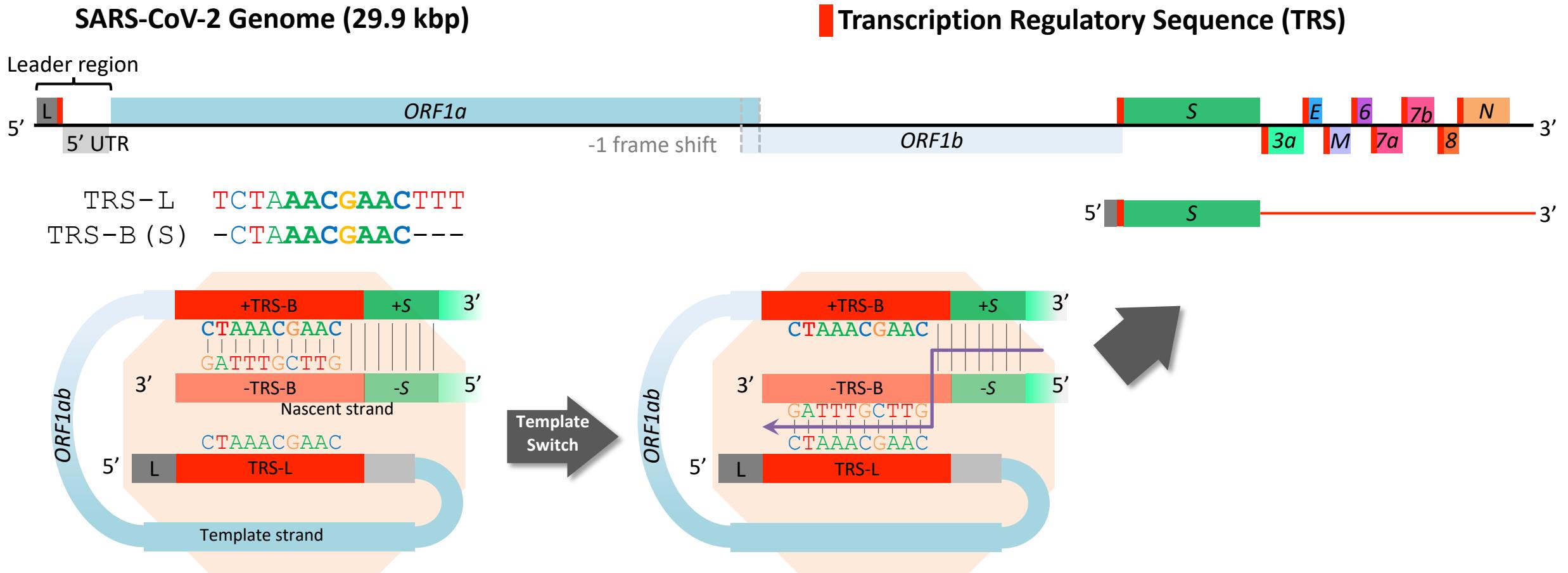
Motivation – Background



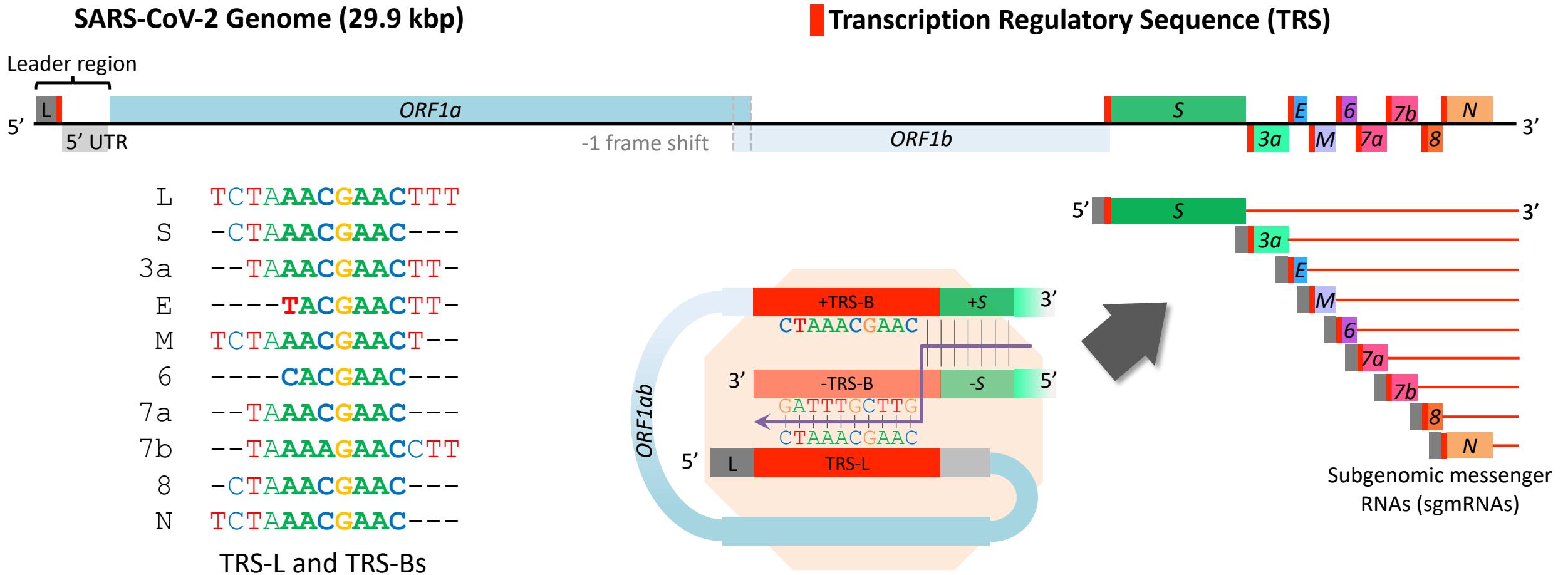
Motivation – Background



Motivation – Background



Motivation – Background



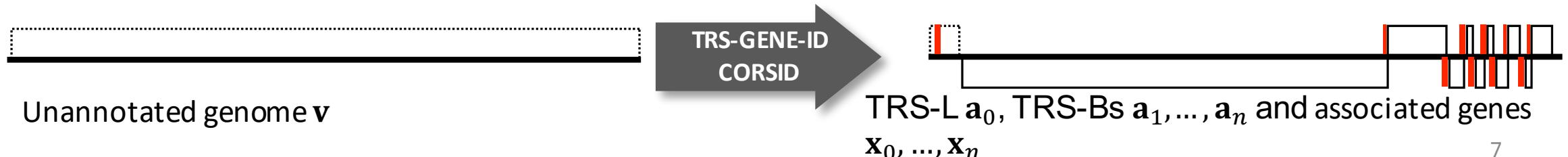
Discontinuous transcription due to template switching of RdRP at transcription regulatory sequences

Motivation – Two Key Questions

Question 1. Can we identify TRS-L and TRS-Bs in **annotated** genomes?



Question 2. Can we identify TRS-L, TRS-Bs and their **corresponding genes** in **unannotated** genomes?



Motivation – Current Methods

- **Motif finding (with annotation):**

- MEME [Bailey et al. 2009] – general purpose
- SuPER [Yang et al. 2021] – needs prior knowledge

- **Gene finding (without annotation):**

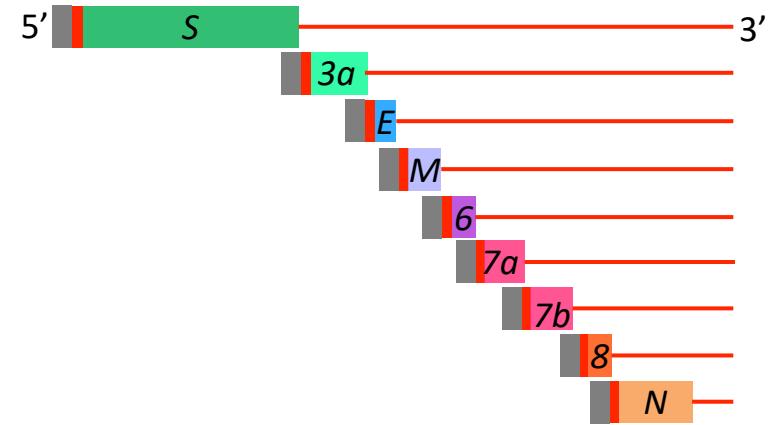
- Prodigal [Hyatt et al. 2010, 2021]
- Glimmer3 [Salzberg et al. 1998, Delcher et al. 2007]
- VADR [Schäffer et al. 2020]

Methods	TRS	Gene
CORSID	✓	✓
CORSID-A	✓	✗
SuPER	✓	✗
MEME	✓	✗
Glimmer3	✗	✓
Prodigal	✗	✓
VADR	✗	✓

There exists no method that combines motif + gene finding for coronaviruses

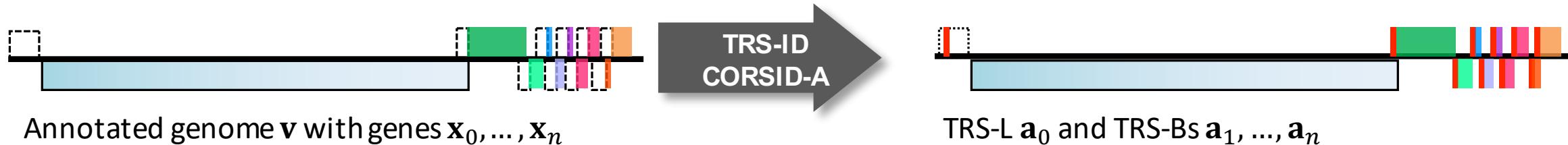
Outline

- TRS Identification
 - Problem Statement
 - Method: CORSID-A
- TRS + Gene Identification
 - Problem Statement
 - Method: CORSID
- Results
- Conclusion



L	TCTAACGAACTTT
S	-CTAACGAAC---
3a	--TAACGAACTT-
E	-----TACGAACTT-
M	TCTAACGAAC T--
6	-----CACGAAC---
7a	---TAACGAAC---
7b	---TAAGAACCTT
8	-CTAACGAAC---
N	TCTAACGAAC---

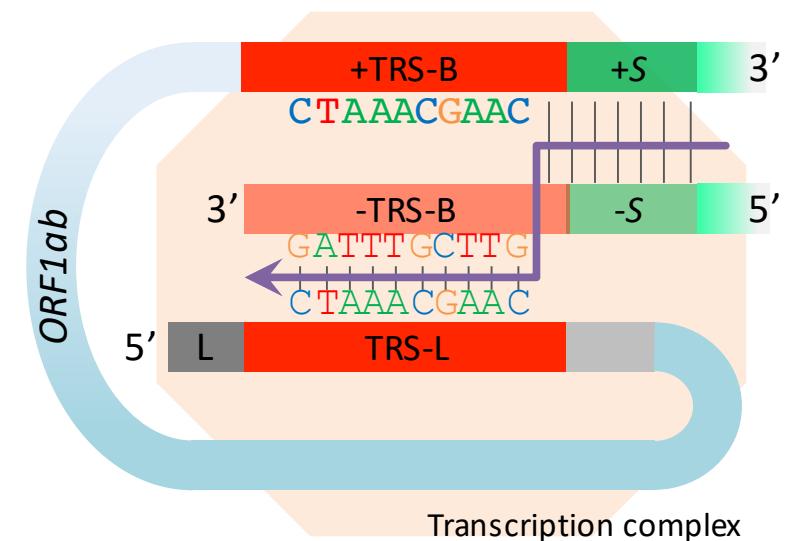
Problem Statement – Definitions



TRS alignment is a constrained multiple sequence alignment (MSA):
(1) No gaps in a_0 ; (2) No internal gaps in a_i ($i > 0$);

$a_0 =$	L	TCTA AACGAAC TTT
$a_1 =$	S	- CTA AACGAAC ---
$a_2 =$	3a	-- TA AACGAAC TT-
$a_3 =$	E	----- TACGAAC TT-
$a_4 =$	M	TCTA AACGAAC T--
$a_5 =$	6	----- CACGAAC ---
$a_6 =$	7a	-- TA AACGAAC ---
$a_7 =$	7b	-- TAAAAGAAC CTT
$a_8 =$	8	- CTAAACGAAC ---
$a_9 =$	N	TCTA AACGAAC ---

Core sequence: $|\mathbf{c}(A)| = 7$

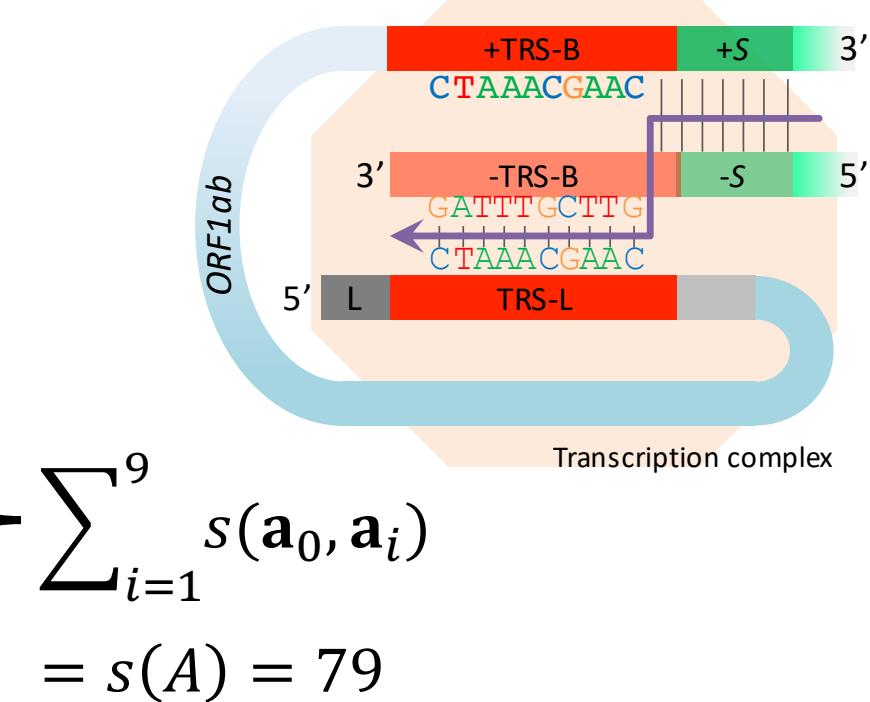


Problem Statement – Definitions

TRS alignment is a constrained multiple sequence alignment (MSA):

- (1) No gaps in \mathbf{a}_0 ; (2) No internal gaps in \mathbf{a}_i ($i > 0$);
- (3) While MSA are typically scored via sum-of-pairs, we score only pairs $(\mathbf{a}_0, \mathbf{a}_i)$

\mathbf{a}_0	L	TRS-L TCTAACGAAC TTT TRS-B	#matches (+1)	#mismatches (-2)	$s(\mathbf{a}_0, \mathbf{a}_i)$
\mathbf{a}_1	S	- CTAACGAAC ---	10	0	10
\mathbf{a}_2	3a	-- TAAACGAAC TT-	11	0	11
\mathbf{a}_3	E	---- TACGAAC TT-	8	1	6
\mathbf{a}_4	M	TCTAACGAAC T--	12	0	12
\mathbf{a}_5	6	---- CACGAAC ---	6	1	4
\mathbf{a}_6	7a	-- TAAACGAAC ---	9	0	9
\mathbf{a}_7	7b	-- TAAAAGAAC CTT	10	2	6
\mathbf{a}_8	8	- CTAACGAAC ---	10	0	10
\mathbf{a}_9	N	TCTAACGAAC ---	11	0	11

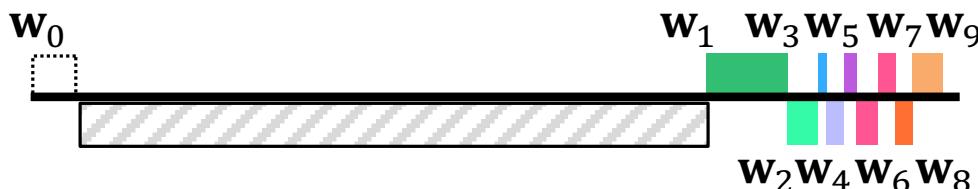


Problem Statement – TRS Identification

TRS Identification Problem: Given non-overlapping sequences w_0, \dots, w_n , core sequence length $\omega > 0$, find a TRS alignment $A = [a_0, \dots, a_n]^T$ such that (i) a_i corresponds to a subsequence in w_i for all $i \in \{0, \dots, n\}$, (ii) the core sequence $c(A)$ has length at least ω and (iii) the alignment has maximum score $s(A)$

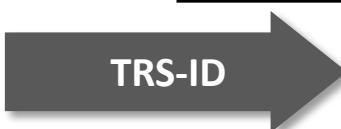
$$s(A) = 79$$
$$c(A) = \text{AACGAAC}$$

Candidate regions: w_0, \dots, w_9



Core sequence $\omega = 7$

$$\begin{aligned} a_0 &= L \quad \text{TCTA} \color{red}{\text{AAC}} \color{green}{\text{GA}} \color{blue}{\text{AC}} \text{TTT} \\ a_1 &= S \quad - \color{blue}{\text{CTA}} \color{red}{\text{AA}} \color{green}{\text{AC}} \color{blue}{\text{GA}} \color{red}{\text{AC}} \text{---} \\ a_2 &= 3a \quad -- \color{red}{\text{TAA}} \color{green}{\text{AC}} \color{blue}{\text{GA}} \color{red}{\text{AC}} \text{TT} - \\ a_3 &= E \quad --- \color{red}{\text{TAC}} \color{green}{\text{GA}} \color{blue}{\text{AC}} \text{TT} - \\ a_4 &= M \quad \text{TCTA} \color{red}{\text{AAC}} \color{green}{\text{GA}} \color{blue}{\text{AC}} \text{T} -- \\ a_5 &= 6 \quad --- \color{blue}{\text{CAC}} \color{green}{\text{GA}} \color{red}{\text{AC}} \text{---} \\ a_6 &= 7a \quad -- \color{red}{\text{TAA}} \color{green}{\text{AC}} \color{blue}{\text{GA}} \text{AC} \text{---} \\ a_7 &= 7b \quad -- \color{red}{\text{TAA}} \color{green}{\text{AA}} \color{blue}{\text{GA}} \color{red}{\text{AC}} \text{TT} \\ a_8 &= 8 \quad - \color{blue}{\text{CTA}} \color{red}{\text{AAC}} \color{green}{\text{GA}} \color{blue}{\text{AC}} \text{---} \\ a_9 &= N \quad \text{TCTA} \color{red}{\text{AAC}} \color{green}{\text{GA}} \color{blue}{\text{AC}} \text{---} \end{aligned}$$

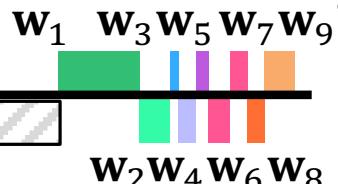


Problem Statement – TRS Identification

TRS Identification Problem: Given non-overlapping sequences w_0, \dots, w_n , core sequence length $\omega > 0$, find a TRS alignment $A = [a_0, \dots, a_n]^T$ such that (i) a_i corresponds to a subsequence in w_i for all $i \in \{0, \dots, n\}$, (ii) the core sequence $c(A)$ has length at least ω and (iii) the alignment has maximum score $s(A)$

$$s(A) = 79$$
$$c(A) = \text{AACGAAC}$$

Candidate regions: w_0, \dots, w_9



Core sequence $\omega = 7$	
$a_0 = L$	TCTA AACGAAC TTT
$a_1 = S$	-CTA AACGAAC -
$a_2 = 3a$	--TA AACGAAC TT-
$a_3 = E$	---TAC GAAC TT-
$a_4 = M$	TCTA AACGAAC T--
$a_5 = 6$	---CAC GAAC ---
$a_6 = 7a$	--TA AACGAAC ---
$a_7 = 7b$	--TA AAAAGAAC CTT
$a_8 = 8$	-CTA AACGAAC ---
$a_9 = N$	TCTA AACGAAC ---



Key observation:
TRS alignment A must have induced core sequence $c(A)$ with length at least ω . Thus, input sequences w_1, \dots, w_n depend on one another and we cannot consider them in isolation.



Method – TRS Identification (I)

Key idea 1:

Break dependency using a subsequence \mathbf{u} of \mathbf{w}_0 requiring that $\mathbf{c}(A)$ contains \mathbf{u} .

Key idea 2:

Solution to TRS-ID corresponds to window \mathbf{u}^* that induces TRS alignment with maximum score.

Constrained TRS Identification Problem:

Given candidate regions $\mathbf{w}_0, \mathbf{w}_1, \dots \mathbf{w}_n$ and subsequence \mathbf{u} of \mathbf{w}_0 , find the optimal TRS alignment such that \mathbf{u} is contained in the core sequence.

Method – TRS Identification (I)

Key idea 1:

Break dependency using a subsequence \mathbf{u} of \mathbf{w}_0 requiring that $\mathbf{c}(A)$ contains \mathbf{u} .

Key idea 2:

Solution to TRS-ID corresponds to window \mathbf{u}^* that induces TRS alignment with maximum score.

Constrained TRS Identification Problem:

Given candidate regions $\mathbf{w}_0, \mathbf{w}_1, \dots \mathbf{w}_n$ and subsequence \mathbf{u} of \mathbf{w}_0 , find the optimal TRS alignment such that \mathbf{u} is contained in the core sequence.

Key idea 3:

Each problem decomposes into n subproblems of aligning \mathbf{w}_0 to \mathbf{w}_i spanning \mathbf{u} .

Key idea 4:

Pairwise Constrained TRS-ID is a variant of local alignment with three differences, (i) alignment $A = [\mathbf{a}_0, \mathbf{a}_i]^T$ may not contain gaps, (ii) \mathbf{a}_0 must contain \mathbf{u} .

Intermezzo – Global vs. Local Alignment

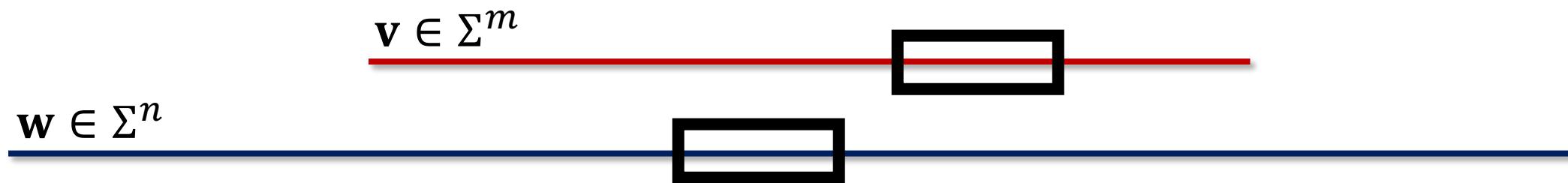
Global Alignment problem: Given strings $v \in \Sigma^m$ and $w \in \Sigma^n$ and scoring function δ , find alignment of v and w with maximum score.

$$\underline{v \in \Sigma^m}$$

$$\underline{w \in \Sigma^n}$$

Intermezzo – Global vs. Local Alignment

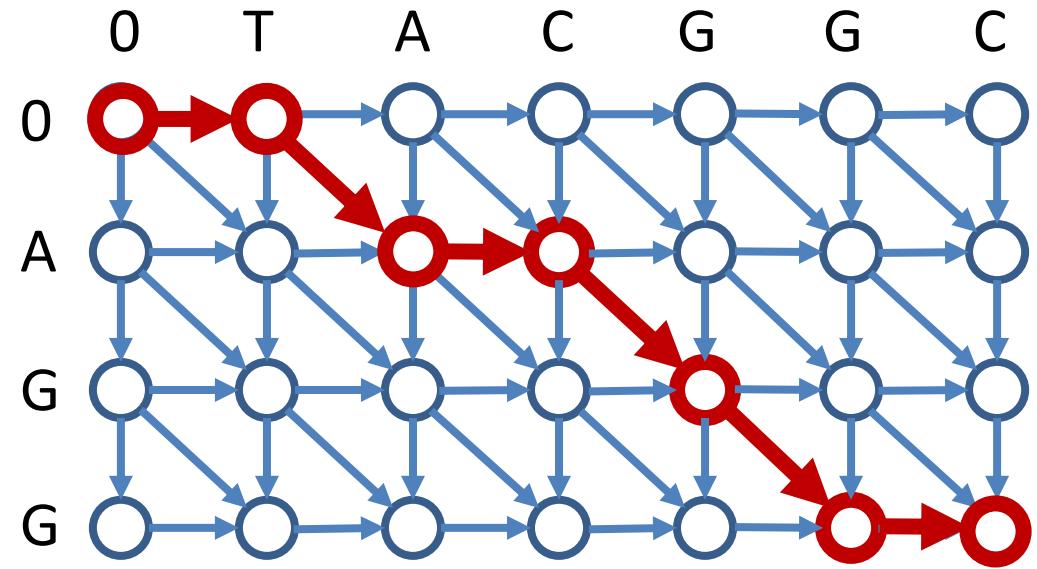
Global Alignment problem: Given strings $v \in \Sigma^m$ and $w \in \Sigma^n$ and scoring function δ , find alignment of v and w with maximum score.



Local Alignment problem: Given strings $v \in \Sigma^m$ and $w \in \Sigma^n$ and scoring function δ , find a substring of v and a substring of w whose alignment has maximum global alignment score s^* among *all* global alignments of *all* substrings of v and w .

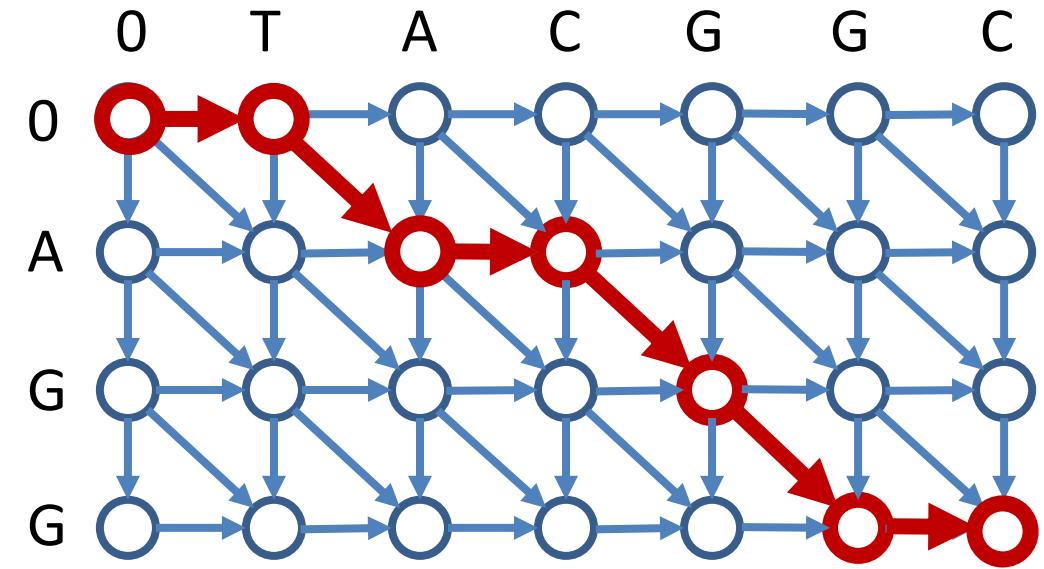
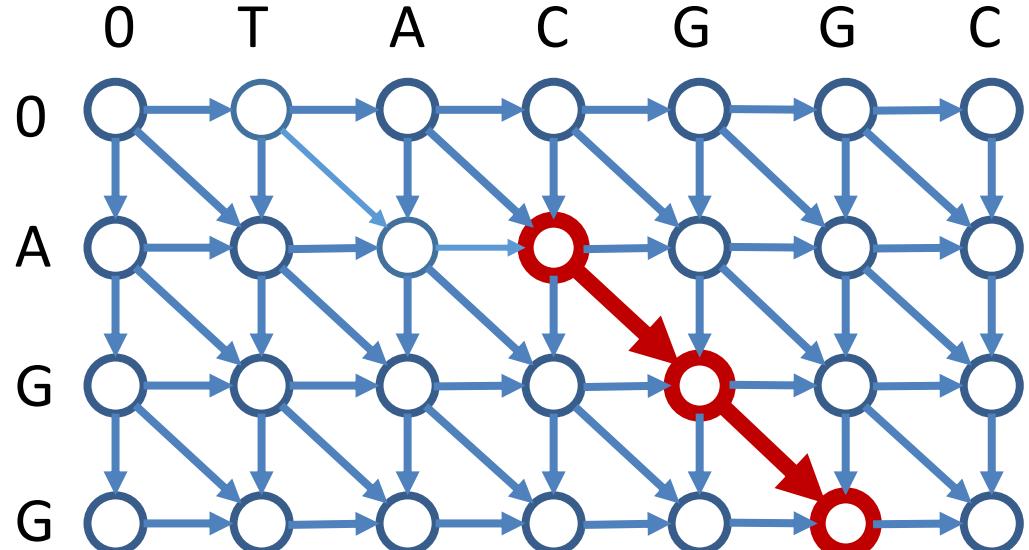
Intermezzo – Global vs. Local Alignment

$$s[i, j] = \max \begin{cases} 0, & \text{if } i = 0 \text{ and } j = 0, \\ s[i - 1, j] + \delta(v_i, -), & \text{if } i > 0, \\ s[i, j - 1] + \delta(-, w_j), & \text{if } j > 0, \\ s[i - 1, j - 1] + \delta(v_i, w_j), & \text{if } i > 0 \text{ and } j > 0. \end{cases}$$



Intermezzo – Global vs. Local Alignment

$$s[i, j] = \max \begin{cases} 0, & \text{if } i = 0 \text{ and } j = 0, \\ s[i - 1, j] + \delta(v_i, -), & \text{if } i > 0, \\ s[i, j - 1] + \delta(-, w_j), & \text{if } j > 0, \\ s[i - 1, j - 1] + \delta(v_i, w_j), & \text{if } i > 0 \text{ and } j > 0. \end{cases}$$



$$s[i, j] = \max \begin{cases} 0, & \text{if } i = 0 \text{ and } j = 0, \\ s[i - 1, j] + \delta(v_i, -), & \text{if } i > 0, \\ s[i, j - 1] + \delta(-, w_j), & \text{if } j > 0, \\ s[i - 1, j - 1] + \delta(v_i, w_j), & \text{if } i > 0 \text{ and } j > 0. \end{cases}$$

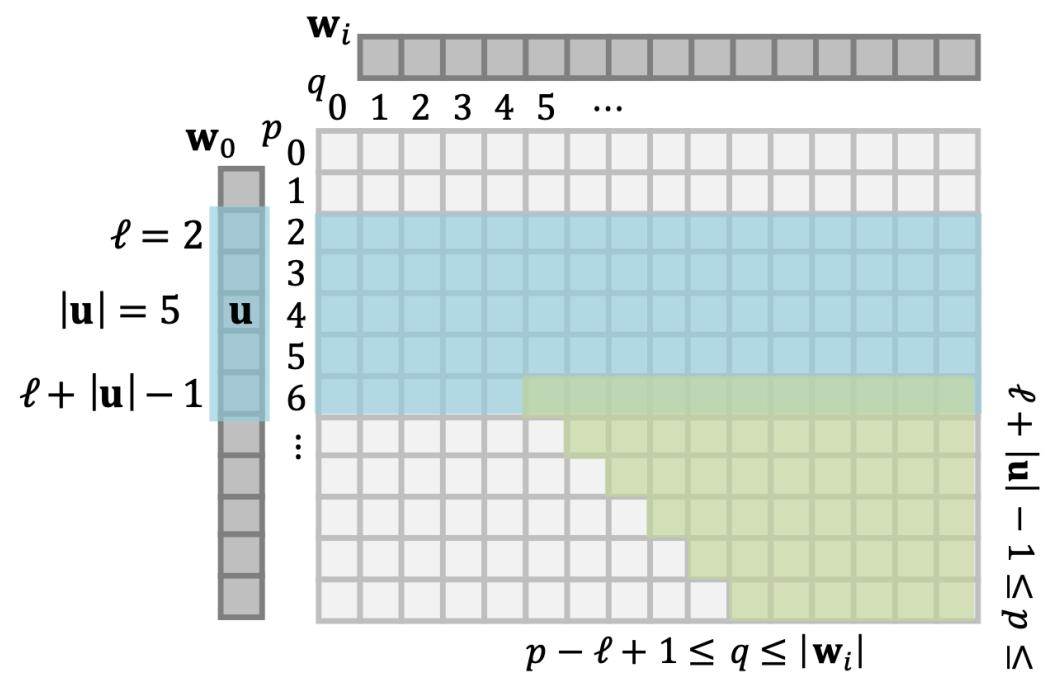
$$s^* = \max_{i,j} s[i, j]$$

Method – TRS Identification (II)

Key idea 4:

Pairwise Constrained TRS-ID is a variant of local alignment with three differences, (i) alignment $A = [\mathbf{a}_0, \mathbf{a}_i]^T$ may not contain gaps, (ii) \mathbf{a}_0 must contain $\mathbf{u} = w_{0,\ell}, \dots, w_{0,\ell+\omega-1}$ at some position ℓ .

$s[p, q]$ is the optimal score of aligning
 $w_{0,1}, \dots, w_{0,p}$ to $w_{i,1}, \dots, w_{i,q}$



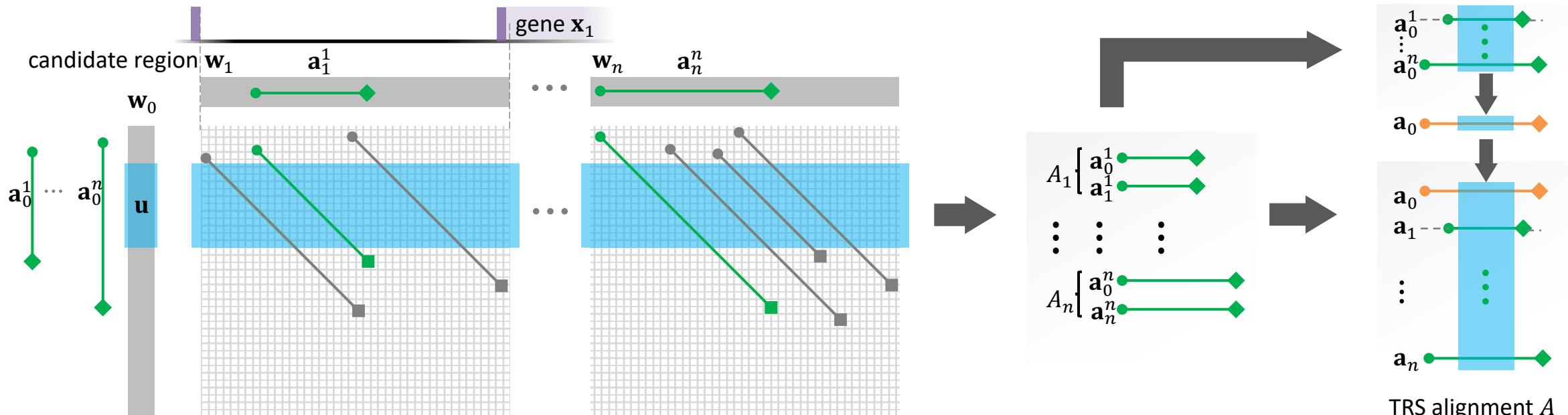
$$s[p, q] = \begin{cases} 0, & \text{if } p = 0 \text{ or } q = 0, \\ \max \{0, s[p - 1, q - 1] + \delta(w_{0,p}, w_{i,q})\}, & \text{if } 1 \leq p < \ell \text{ and } q \geq 1, \\ s[p - 1, q - 1] + \delta(w_{0,p}, w_{i,q}), & \text{if } p \geq \ell \text{ and } q \geq 1, \end{cases} \quad (2)$$

$$(p^*, q^*) = \arg \max_{\ell+|\mathbf{u}|-1 \leq p \leq |\mathbf{w}_0|, p-\ell+1 \leq q \leq |\mathbf{w}_i|} s[p, q]. \quad (3)$$

Methods – CORSID-A

Constrained TRS Identification Problem:

Given candidate regions w_0, w_1, \dots, w_n and subsequence u of w_0 , find the optimal TRS alignment such that u is contained in the core sequence.

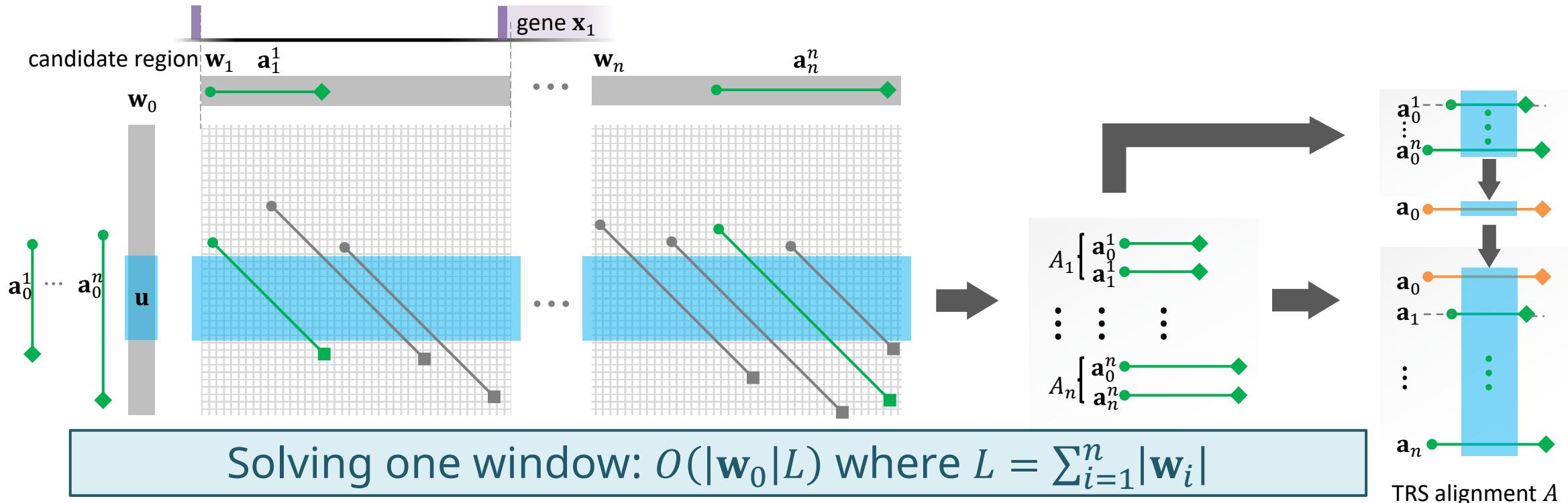


Solving one window: $O(|w_0|L)$ where $L = \sum_{i=1}^n |w_i|$

Methods – CORSID-A

Constrained TRS Identification Problem:

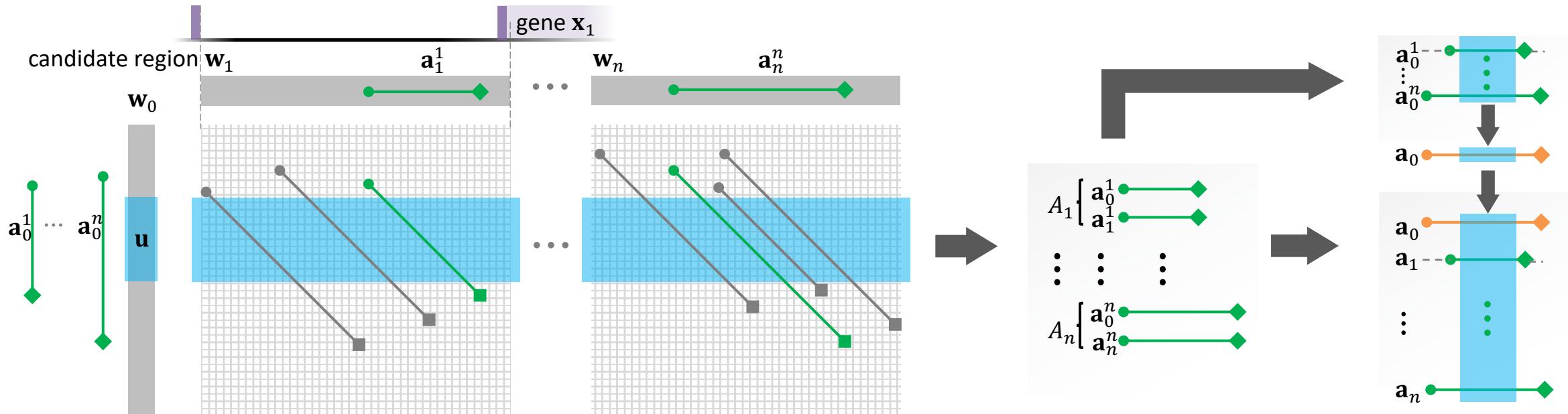
Given candidate regions w_0, w_1, \dots, w_n and subsequence u of w_0 , find the optimal TRS alignment such that u is contained in the core sequence.



Methods – CORSID-A

Constrained TRS Identification Problem:

Given candidate regions w_0, w_1, \dots, w_n and subsequence u of w_0 , find the optimal TRS alignment such that u is contained in the core sequence.

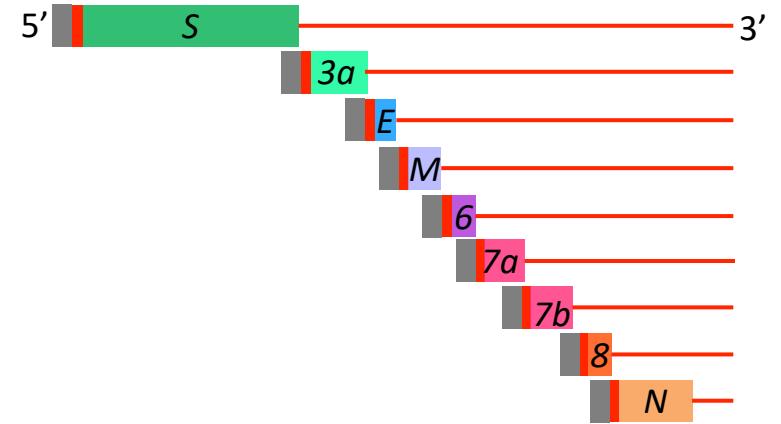


Solving one window: $O(|w_0|L)$ where $L = \sum_{i=1}^n |w_i|$

Solving all $O(|w_0|)$ sliding windows: $O(|w_0|^2 L)$ where $L = \sum_{i=1}^n |w_i|$

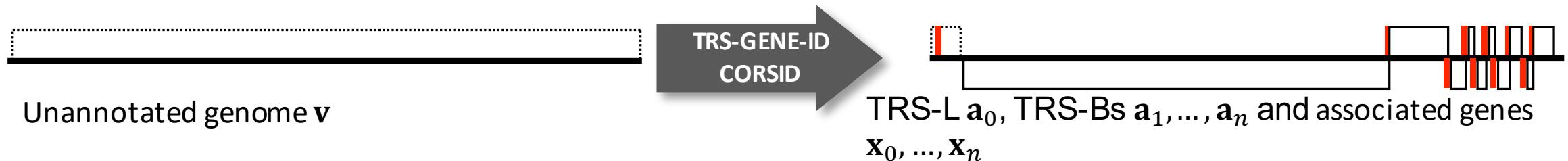
Outline

- TRS Identification
 - Problem Statement
 - Method: CORSID-A
- TRS + Gene Identification
 - Problem Statement
 - Method: CORSID
- Results
- Conclusion

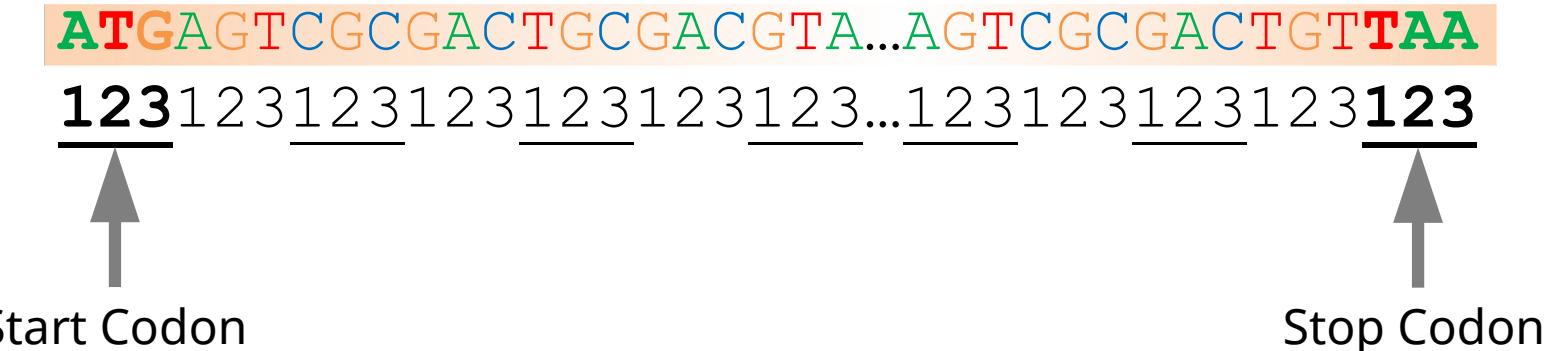


L	TCTAACGAACTTT
S	-CTAACGAAC---
3a	--TAACGAACTT-
E	-----TACGAACTT-
M	TCTAACGAAC T--
6	-----CACGAAC---
7a	---TAACGAAC---
7b	---TAAGAACCTT
8	-CTAACGAAC---
N	TCTAACGAAC---

Problem Statement – TRS-Gene Identification



- Enumerate ORFs for genes



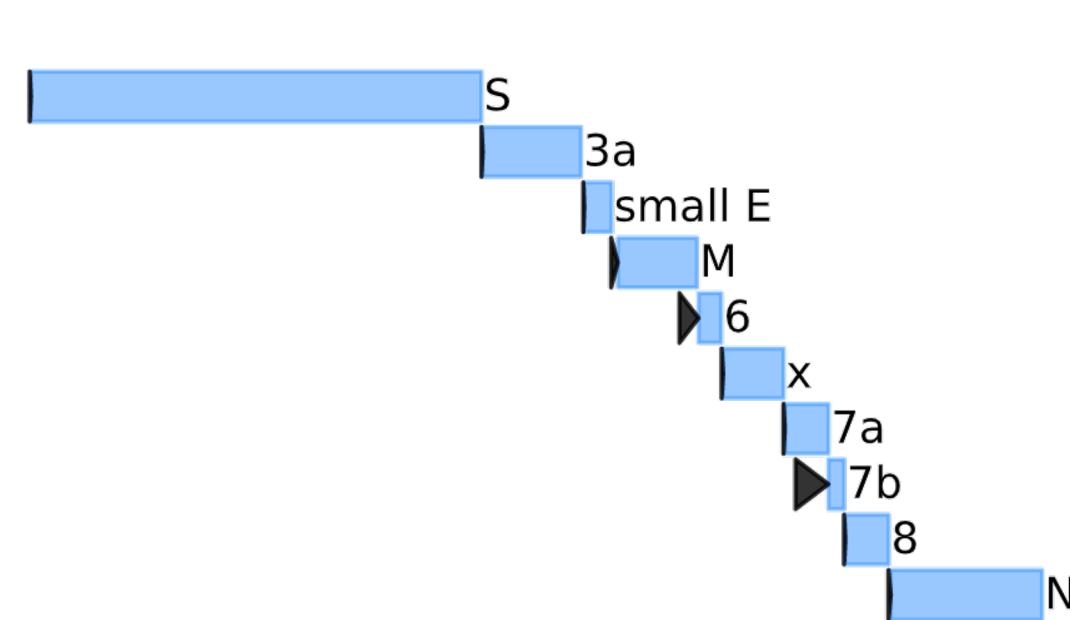
- Genes should not overlap – Independent set problem
- Genome coverage $g(A)$ should be large
- Genes should have similar TRS sequences

Problem Statement – TRS-Gene Identification



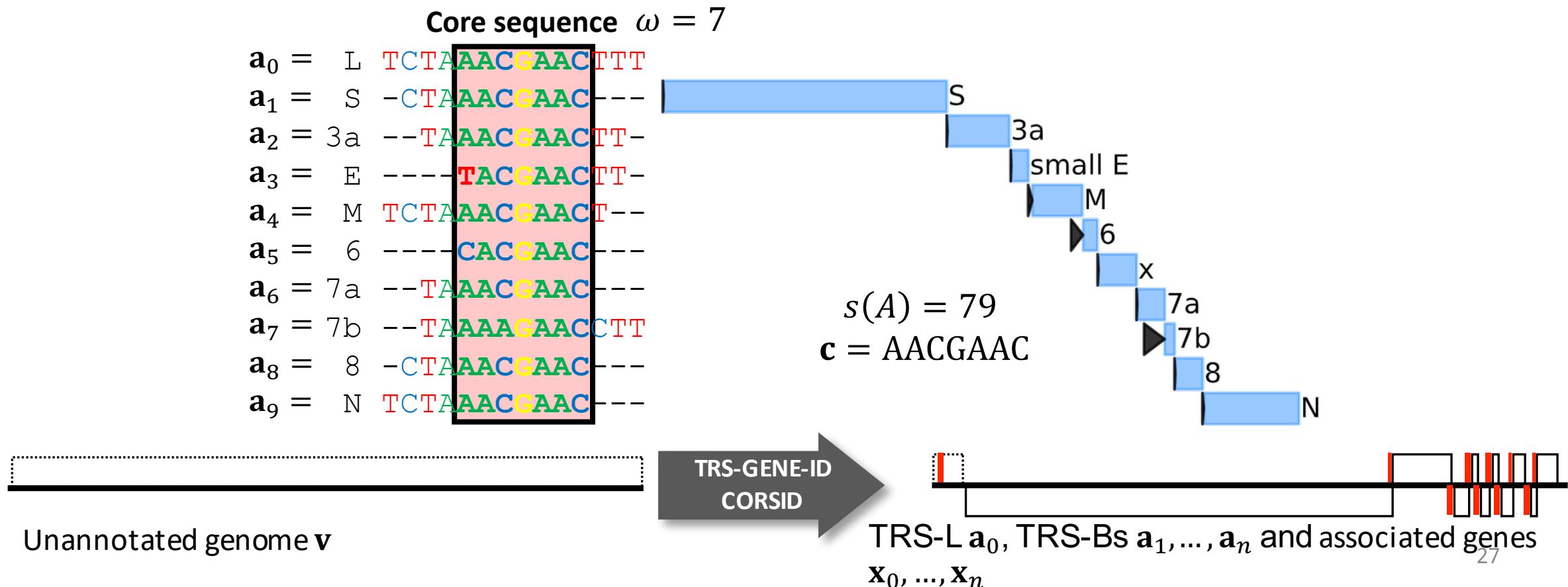
- TRS alignment A induces a gene set $\Gamma(A)$

$a_0 =$	L	TCTA AACGAAC TTT
$a_1 =$	S	-CTA AACGAAC ---
$a_2 =$	3a	---TA AACGAAC TT-
$a_3 =$	E	-----TACGAAC TT-
$a_4 =$	M	TCTA AACGAAC T---
$a_5 =$	6	-----CACGAAC ---
$a_6 =$	7a	---TA AACGAAC ---
$a_7 =$	7b	---TA AAAGAAC CTT
$a_8 =$	8	-CTA AACGAAC ---
$a_9 =$	N	TCTA AACGAAC ---



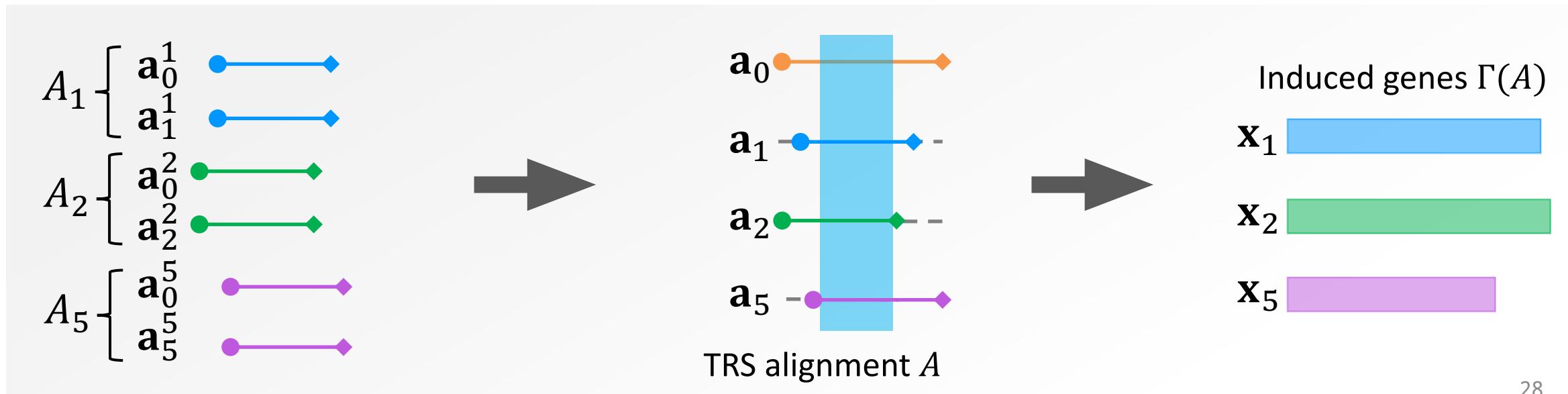
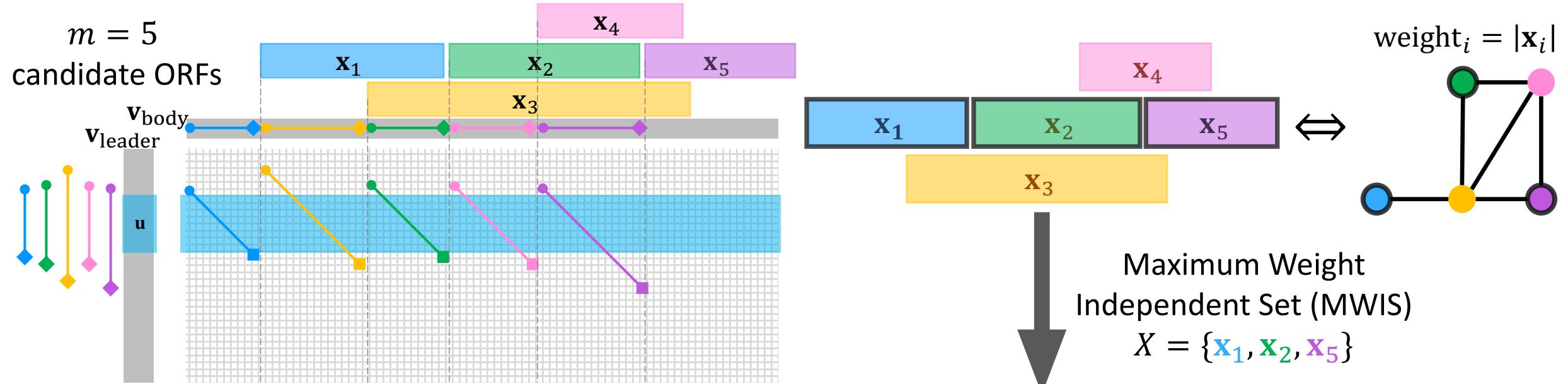
Problem Statement – TRS-Gene Identification

TRS-Gene Identification Problem: Given leader region v_{leader} , body region v_{body} , core-sequence length $\omega > 0$, find a TRS alignment $A = [a_i]$ such that (i) a_0 corresponds to a subsequence in v_{leader} , (ii) a_i corresponds to a subsequence in v_{body} for all $i \geq 1$, (iii) the core sequence $c(A)$ has length at least ω , (iv) A induces a set $\Gamma(A)$ of genes with maximum genome coverage $g(A)$ and subsequently maximum score $s(A)$.



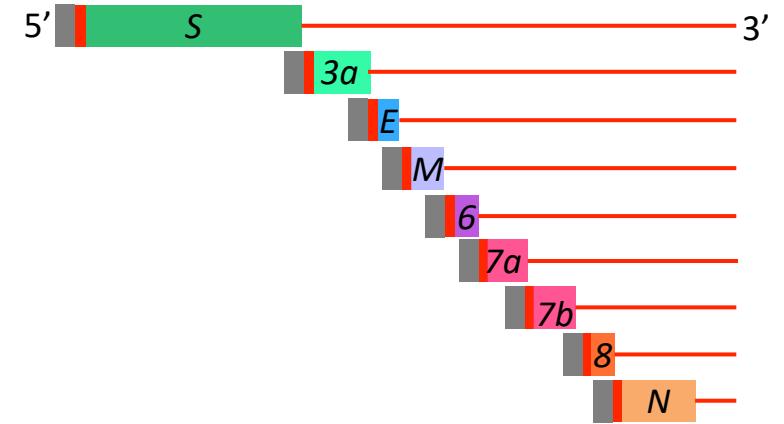
Methods – CORSID

$$O(|\mathbf{v}_{\text{leader}}|^2 |\mathbf{v}_{\text{body}}| + |\mathbf{v}_{\text{leader}}| \cdot m)$$



Outline

- TRS Identification
 - Problem Statement
 - Method: CORSID-A
- TRS + Gene Identification
 - Problem Statement
 - Method: CORSID
- Results
- Conclusion

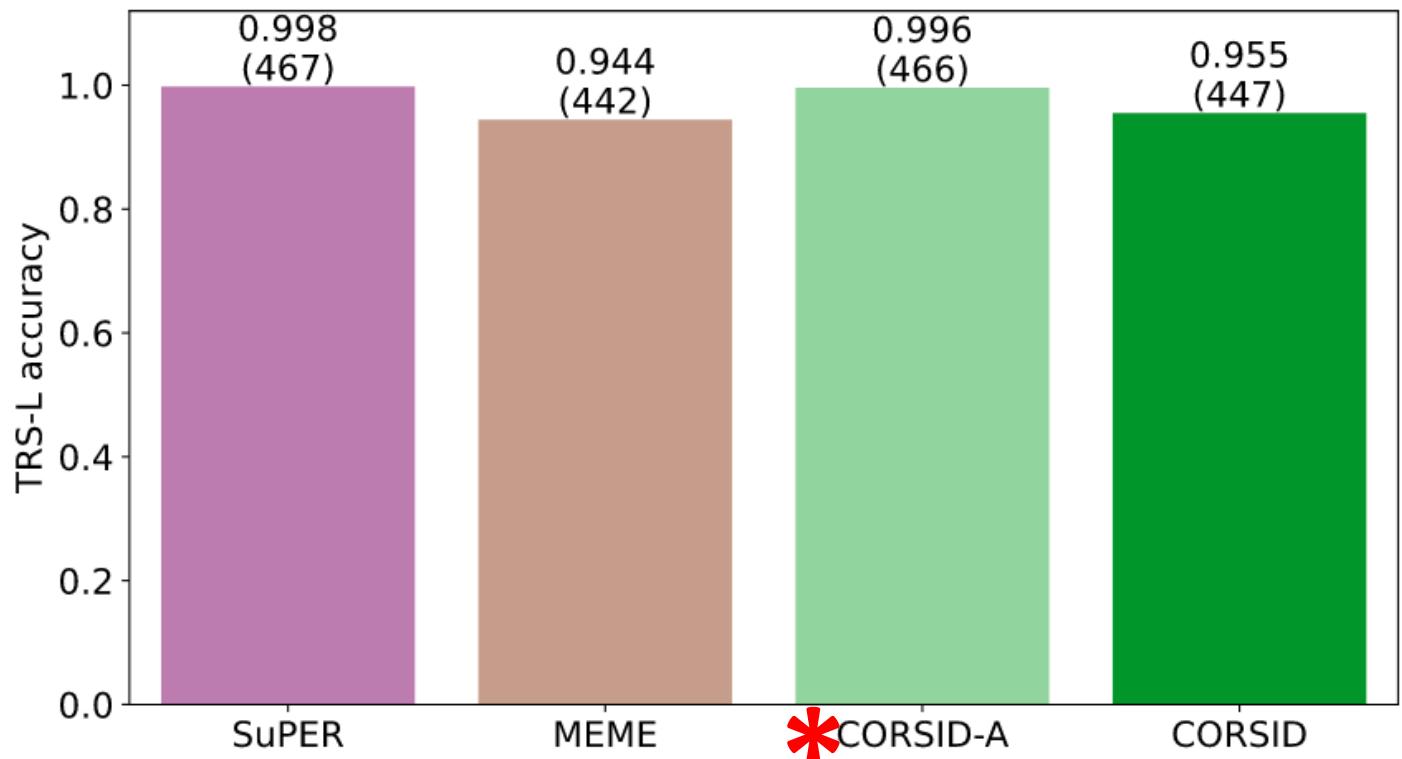


L	TCTAACGAACTTT
S	-CTAACGAAC---
3a	--TAACGAACTT-
E	-----TACGAACTT-
M	TCTAACGAAC T--
6	-----CACGAAC---
7a	--TAACGAAC---
7b	--TAAAAGAACCTT
8	-CTAACGAAC---
N	TCTAACGAAC---

Results – TRS-L Accuracy

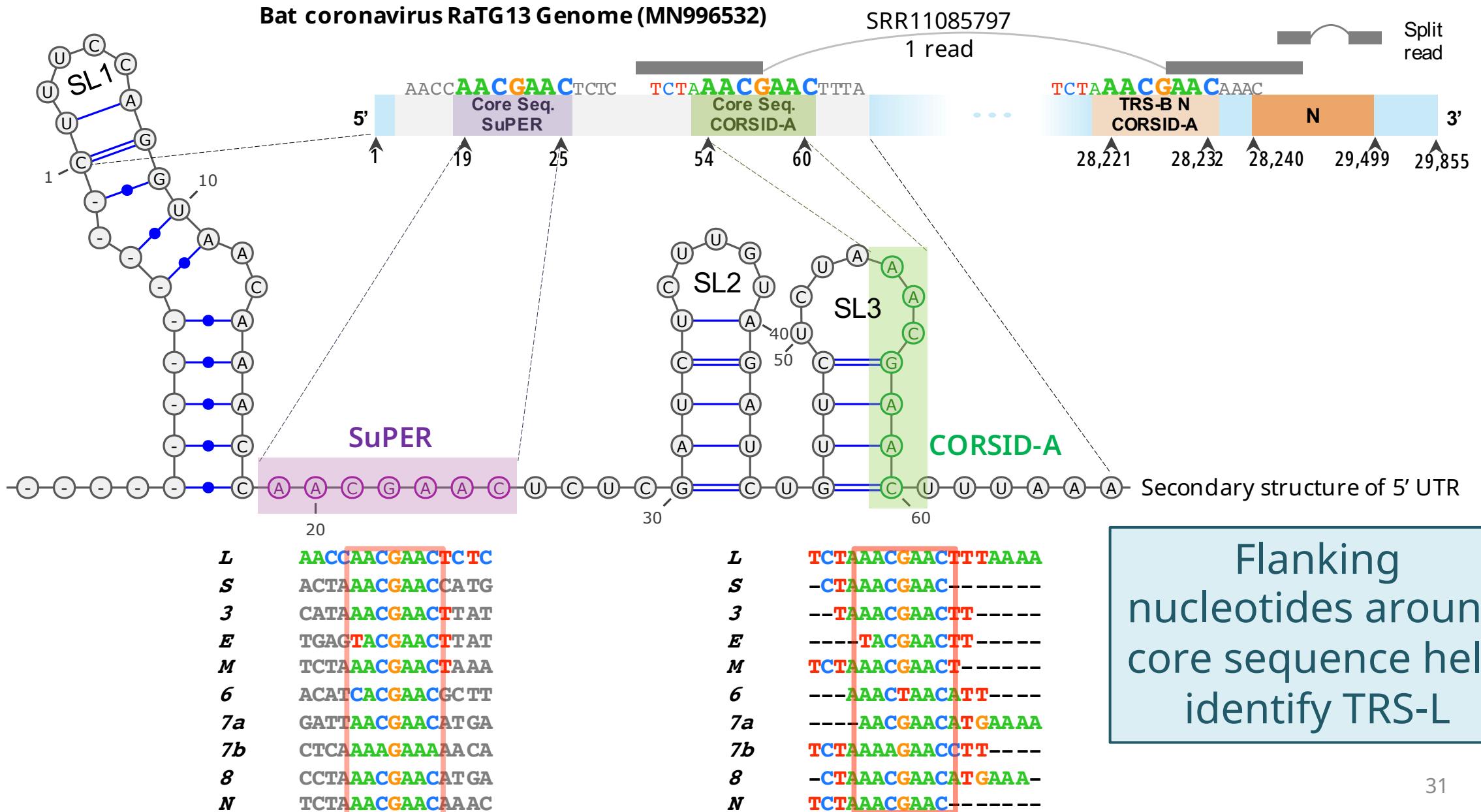
- GenBank: 468 genomes of coronaviruses
- Ground truth based on secondary structure + MSA of leader region

Methods	TRS	Gene	Additional Signal
CORSID	✓	✓	-
CORSID-A	✓	✗	-
SuPER	✓	✗	Secondary structure Predefined motifs
MEME	✓	✗	-
Glimmer3	✗	✓	
Prodigal	✗	✓	
VADR	✗	✓	



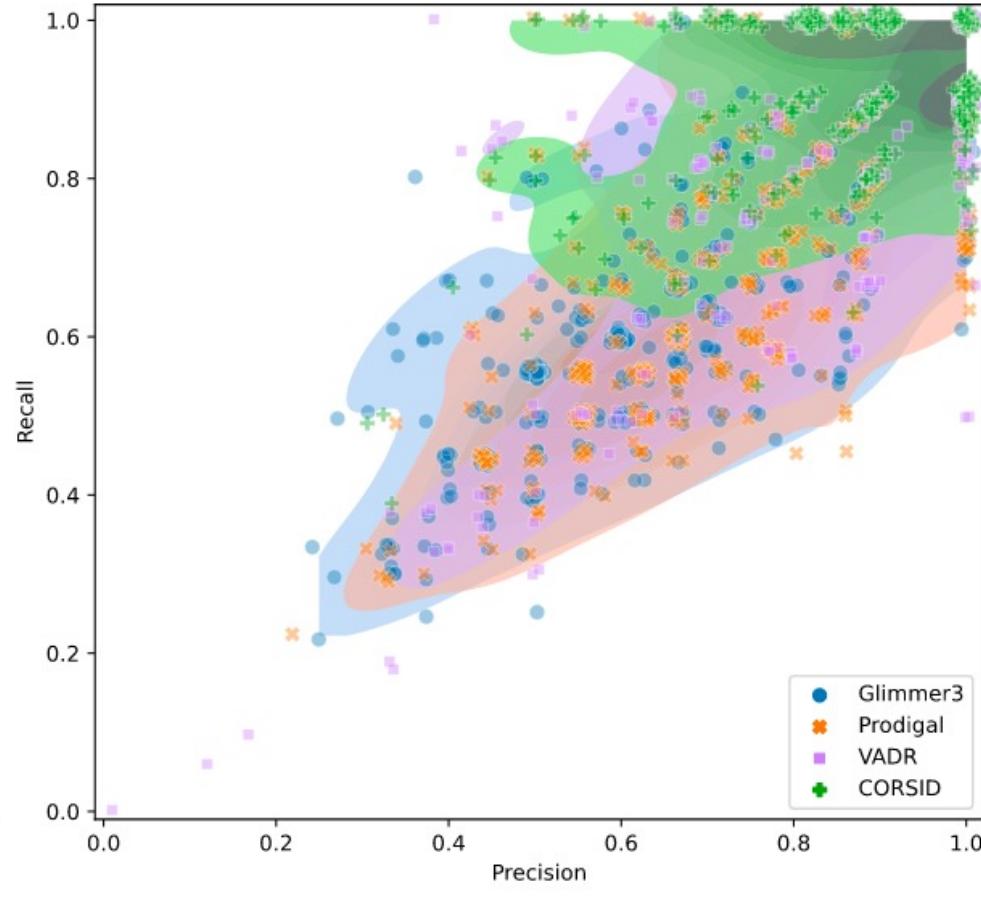
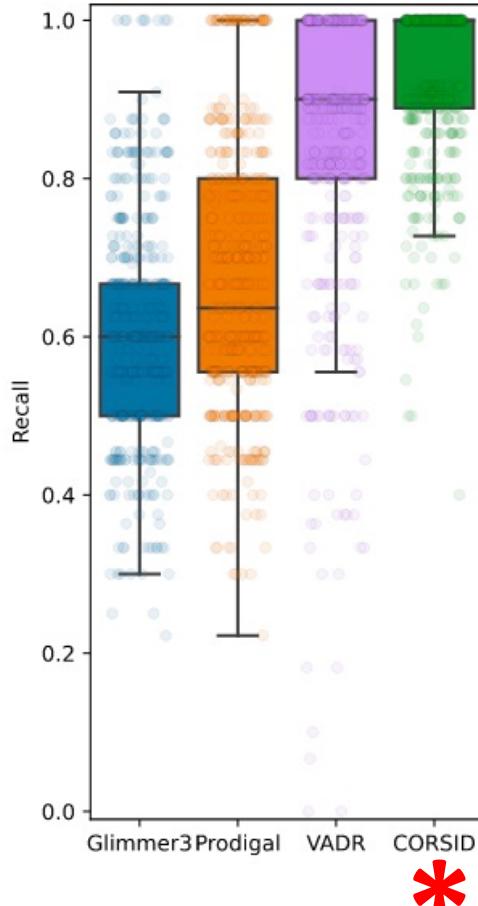
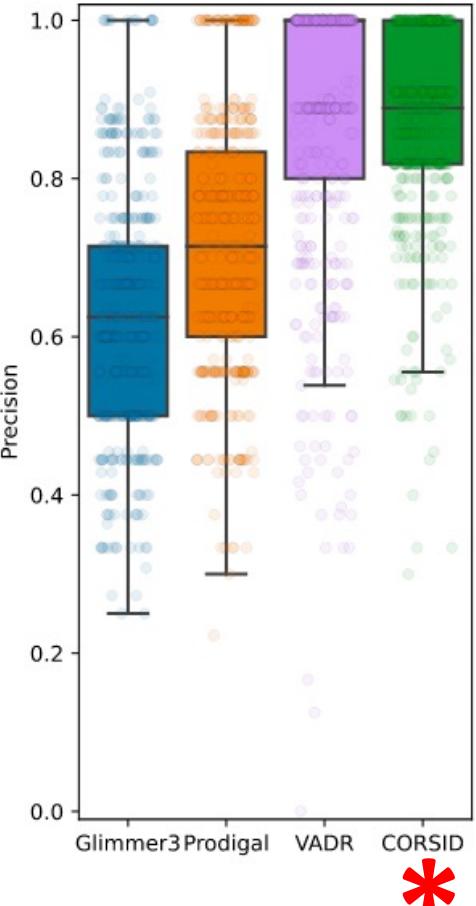
CORSID-A correctly identifies TRS-L

Results – CORSID-A finds the correct TRS-L site



Results – CORSID – Gene Identification

Methods	TRS	Gene
CORSID	✓	✓
CORSID-A	✓	✗
SuPER	✓	✗
MEME	✓	✗
Glimmer3	✗	✓
Prodigal	✗	✓
VADR	✗	✓



Welcome to the CORSID Visualization Tool

In this demo we show 468 samples used in the analysis. You can start exploring results of CORSID/CORSID-A by clicking any link in the table.



Michael Xiang



Yichi Zhang



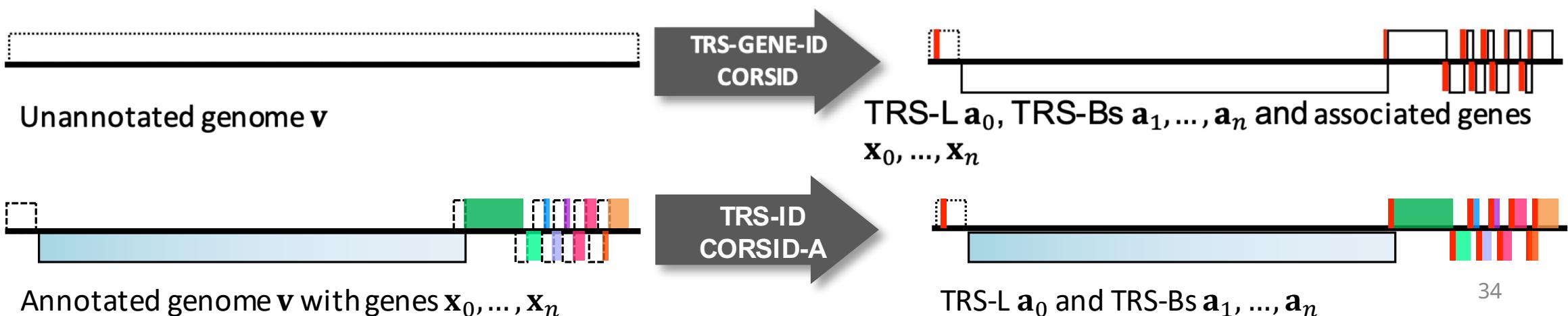
Ayesha Kazi

Search: Genus: Subgenus:

Sample ▲	Genus ▢	Subgenus ▢	CORSID viz	CORSID-A viz
AC_000192	Betacoronavirus	Embecovirus	🔗	🔗
AF220295	Betacoronavirus	Embecovirus	🔗	🔗
AY319651	Gammacoronavirus	Igacovirus	🔗	🔗
AY514485	Gammacoronavirus	Igacovirus	🔗	🔗
DQ011855	Betacoronavirus	Embecovirus	🔗	🔗
DQ022305	Betacoronavirus	Sarbecovirus	🔗	🔗
DQ071615	Betacoronavirus	Sarbecovirus	🔗	🔗
DQ415899	Betacoronavirus	Embecovirus	🔗	🔗
DQ648794	Betacoronavirus	Merbecovirus	🔗	🔗
DQ648856	Betacoronavirus	Sarbecovirus	🔗	🔗
DQ648857	Betacoronavirus	Sarbecovirus	🔗	🔗
DQ811787	Alphacoronavirus	Tegacovirus	🔗	🔗
DQ848678	Alphacoronavirus	Tegacovirus	🔗	🔗
EF065510	Betacoronavirus	Merbecovirus	🔗	🔗
EF065514	Betacoronavirus	Nobecovirus	🔗	🔗
EF065515	Betacoronavirus	Nobecovirus	🔗	🔗
EF065516	Betacoronavirus	Nobecovirus	🔗	🔗
EF446615	Betacoronavirus	Embecovirus	🔗	🔗
EU022525	Gammacoronavirus	Igacovirus	🔗	🔗
EU022526	Gammacoronavirus	Igacovirus	🔗	🔗
EU074218	Alphacoronavirus	Tegacovirus	🔗	🔗
EU186072	Alphacoronavirus	Tegacovirus	🔗	🔗
EU420137	Alphacoronavirus	Minunacovirus	🔗	🔗
EU637854	Gammacoronavirus	Igacovirus	🔗	🔗
EU714029	Gammacoronavirus	Igacovirus	🔗	🔗
EU917407	Gammacoronavirus	Igacovirus	🔗	🔗

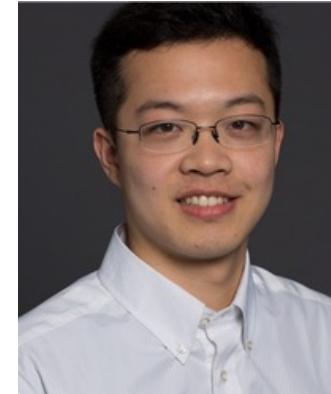
Conclusion

- CORSID:
 - The first method to simultaneously identify TRS sites and genes
 - *De novo* identification of TRS sites and genes
 - Outperforms state-of-the-art methods
- Future direction:
 - Alternative start codon, Kozak sequence
 - RNA-seq, split reads
 - Incomplete genome without TRS-L

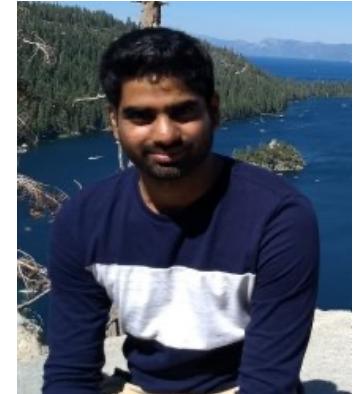


Acknowledgements

- Joint first authors:
Chuanyi Zhang, Palash Sashittal
- Michael Xiang, Yichi Zhang, and Ayesha Kazi
- Grants
 - NSF
 - CRII: CCF-1850502
 - RAPID: CCF-2027669
 - CAREER: CCF-2046488
 - AWS
Greg Gulick Honorary Research Award
- Availability
 - Code
 - WebApp



Chuanyi Zhang



Palash Sashittal



install with bioconda pypi package 0.1.3
<https://github.com/elkebir-group/CORSID>
<https://elkebir-group.github.io/CORSID-viz>