



西安交通大学
XI'AN JIAOTONG UNIVERSITY

AI-powered Structural variation detection

Kai Ye

Xi'an Jiaotong University, China

kaiye@xjtu.edu.cn

Education and Work experience

Education

■ 1995~2003	Wuhan University, CN	Biopharmaceutical Science	B.S.; M.S.
■ 2004~2008	Leiden University, NL	Computer Science	PhD

Work

■ 2008~2009	European Bioinformatics Institute, UK	Supervised by Rolf Apweiler	Postdoc
■ 2009~2012	Leiden University Medical Center, NL	Molecular Epidemiology	Assistant professor
■ 2012~2016	Washington University in St. Louis, US	Genome Institute	Assistant professor
■ 2016~ now	Xi'an Jiaotong University, CN	Automation	Professor

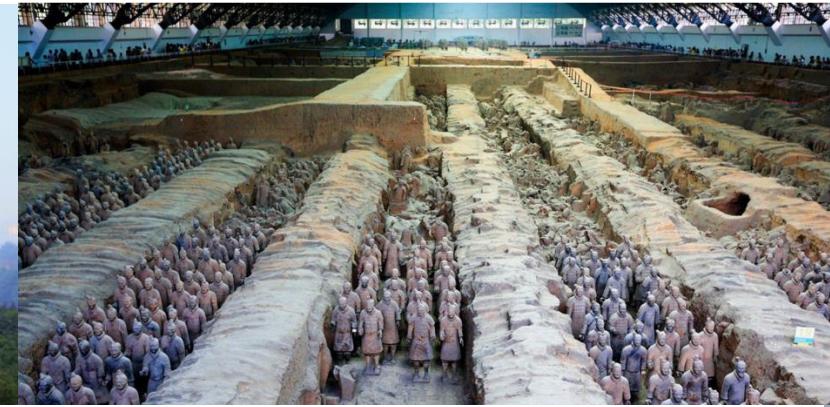
Education and Work experience

Education

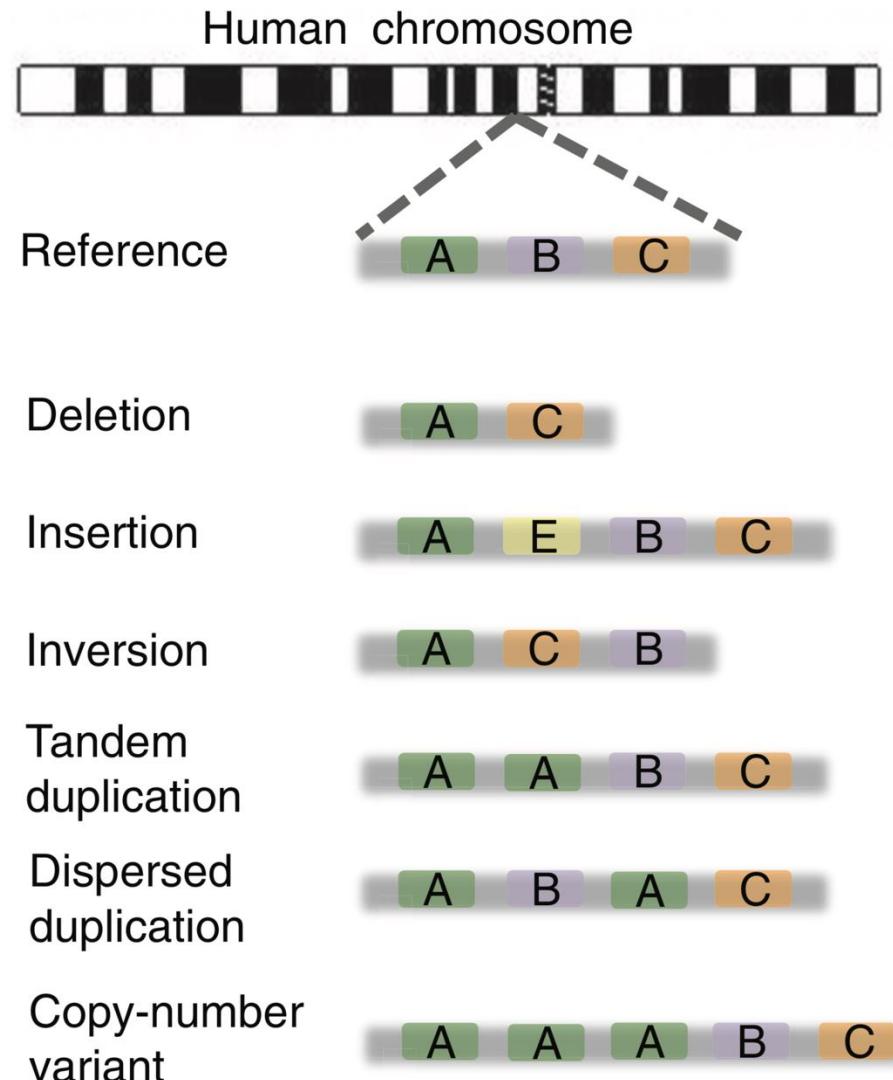
- | | | | |
|-------------|-----------------------|---------------------------|------------|
| ■ 1995~2003 | Wuhan University, CN | Biopharmaceutical Science | B.S.; M.S. |
| ■ 2004~2008 | Leiden University, NL | Computer Science | PhD |

Work

- | | | | |
|-------------|--|-----------------------------|---------------------|
| ■ 2008~2009 | European Bioinformatics Institute, UK | Supervised by Rolf Apweiler | Postdoc |
| ■ 2009~2012 | Leiden University Medical Center, NL | Molecular Epidemiology | Assistant professor |
| ■ 2012~2016 | Washington University in St. Louis, US | Genome Institute | Assistant professor |
| ■ 2016~ now | Xi'an Jiaotong University, CN | Automation | Professor |

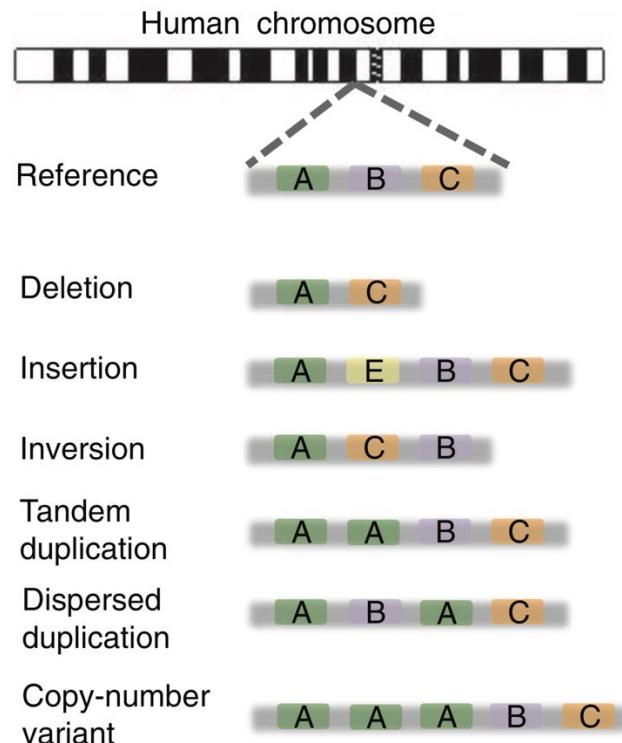


Structural variation (SV)

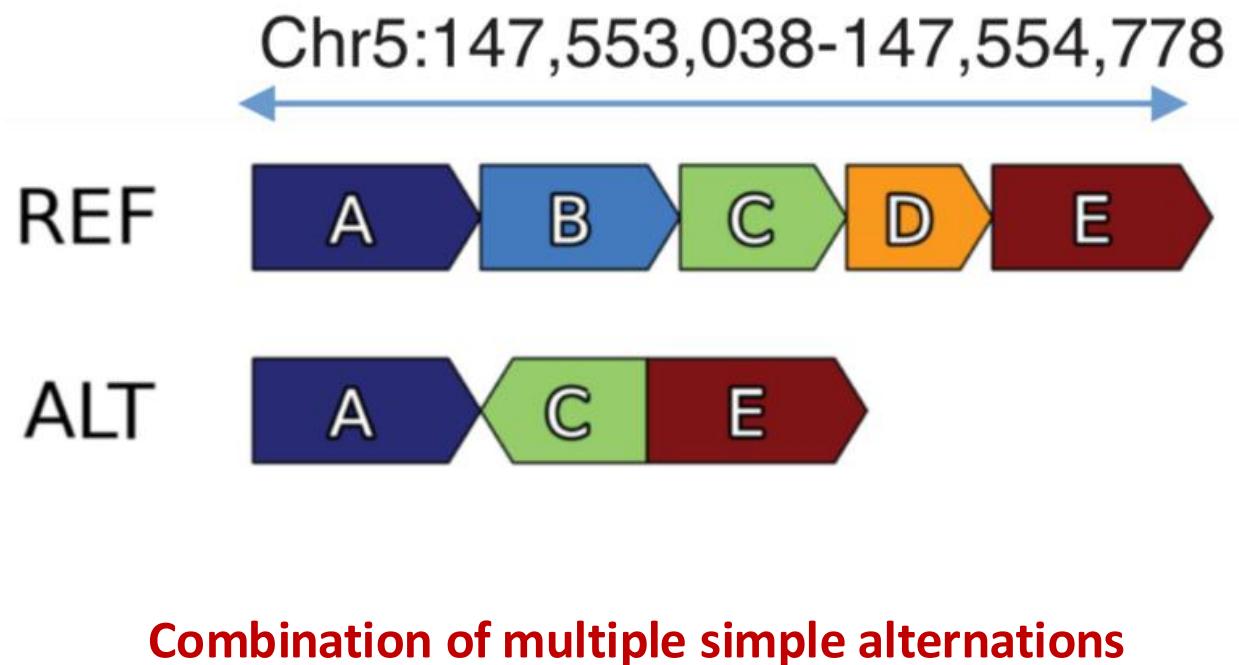


Structural variation (SV)

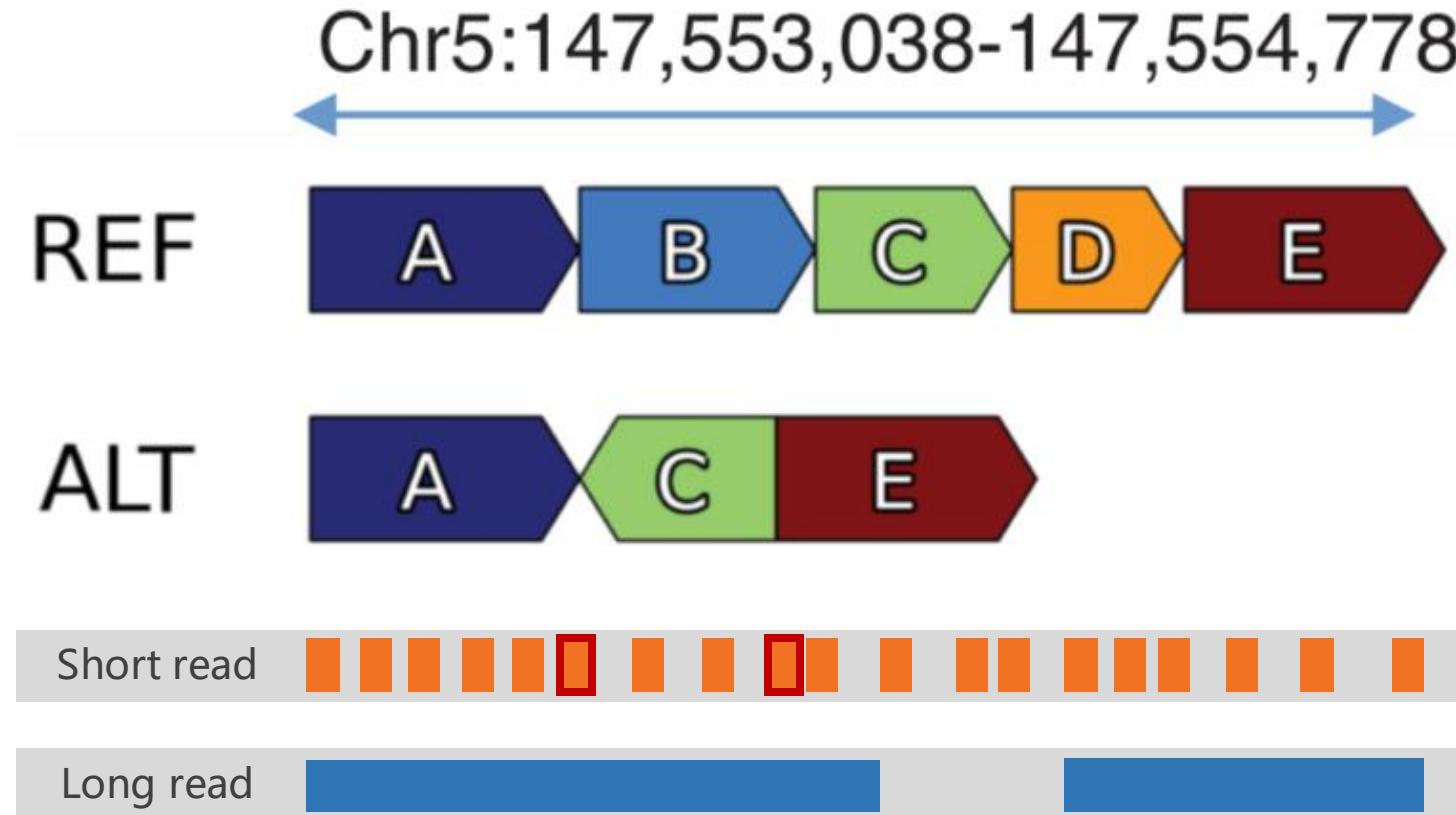
Simple SV



Complex SV



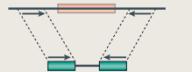
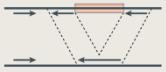
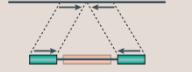
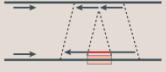
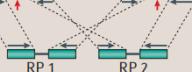
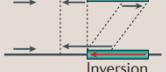
Long reads are better for SV detection



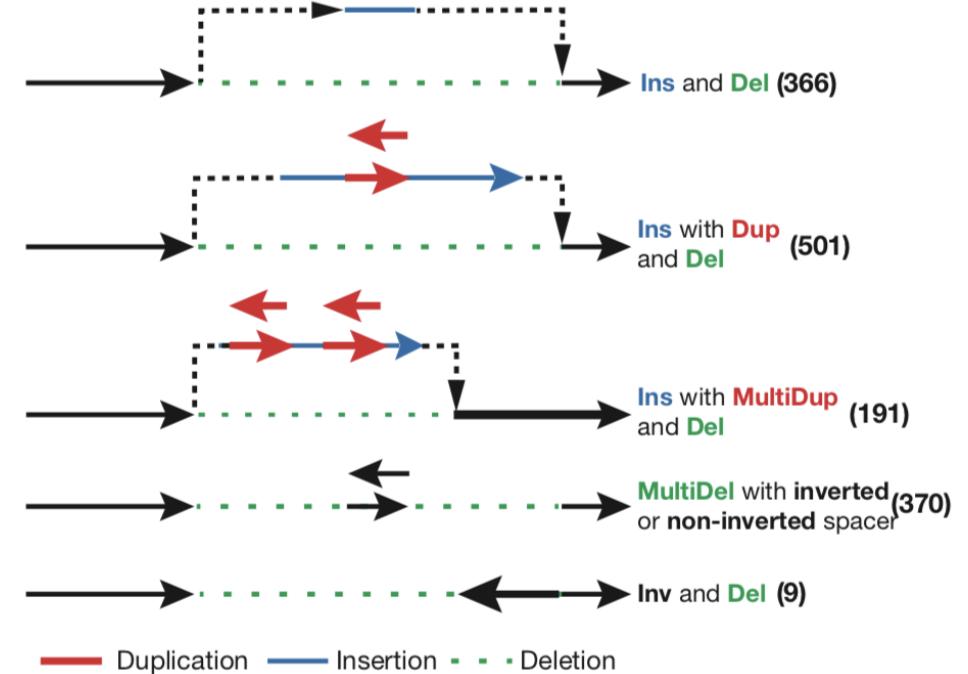
Short reads provide linking information between two adjacent segments
but too short to reveal connection among breakpoints

Challenges of detecting complex SV with long reads (1)

Model based approaches

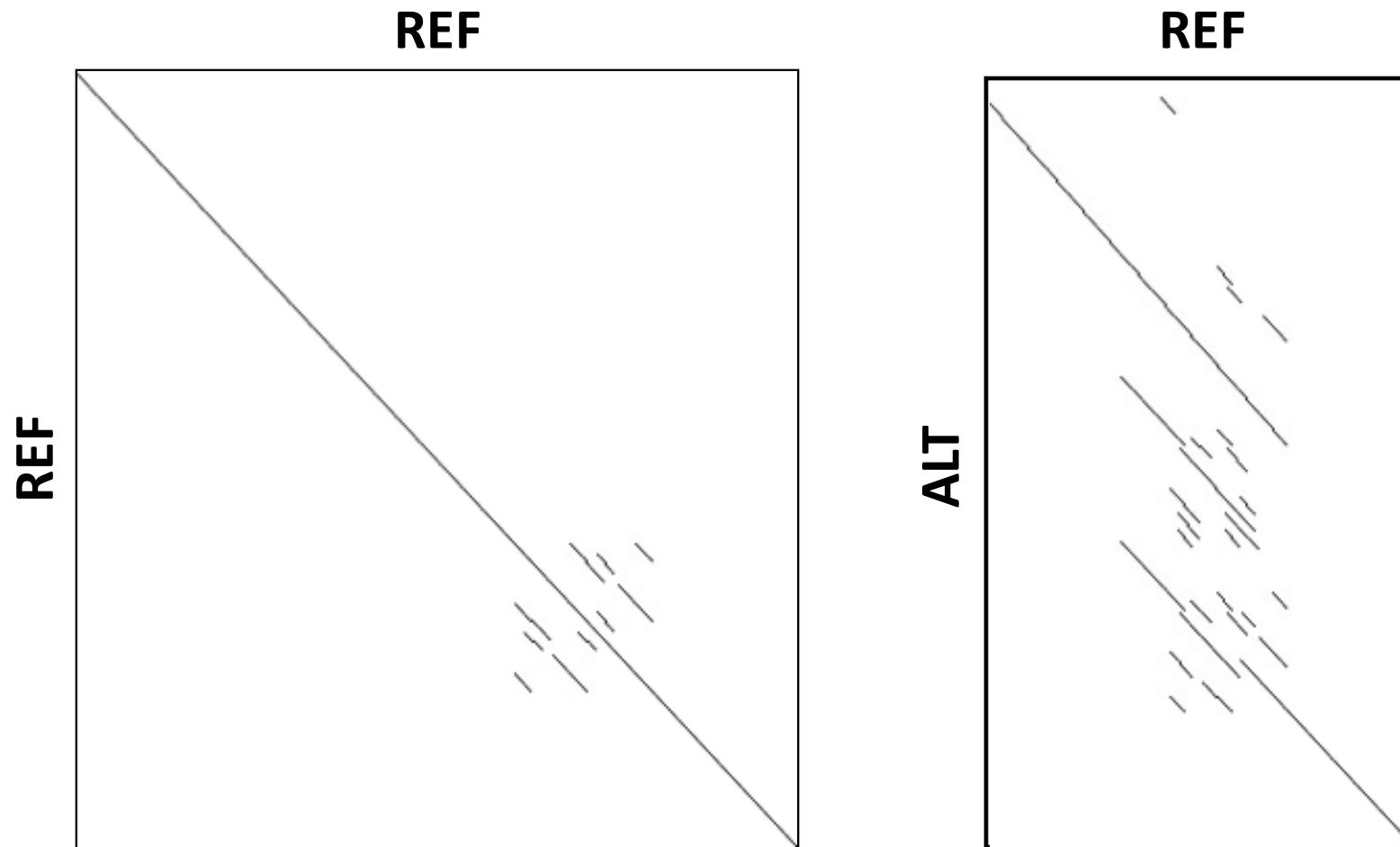
SV classes	Read pair	Read depth	Split read
Deletion			
Novel sequence insertion		Not applicable	
Mobile-element insertion		Not applicable	
Inversion		Not applicable	

Hard to model and code



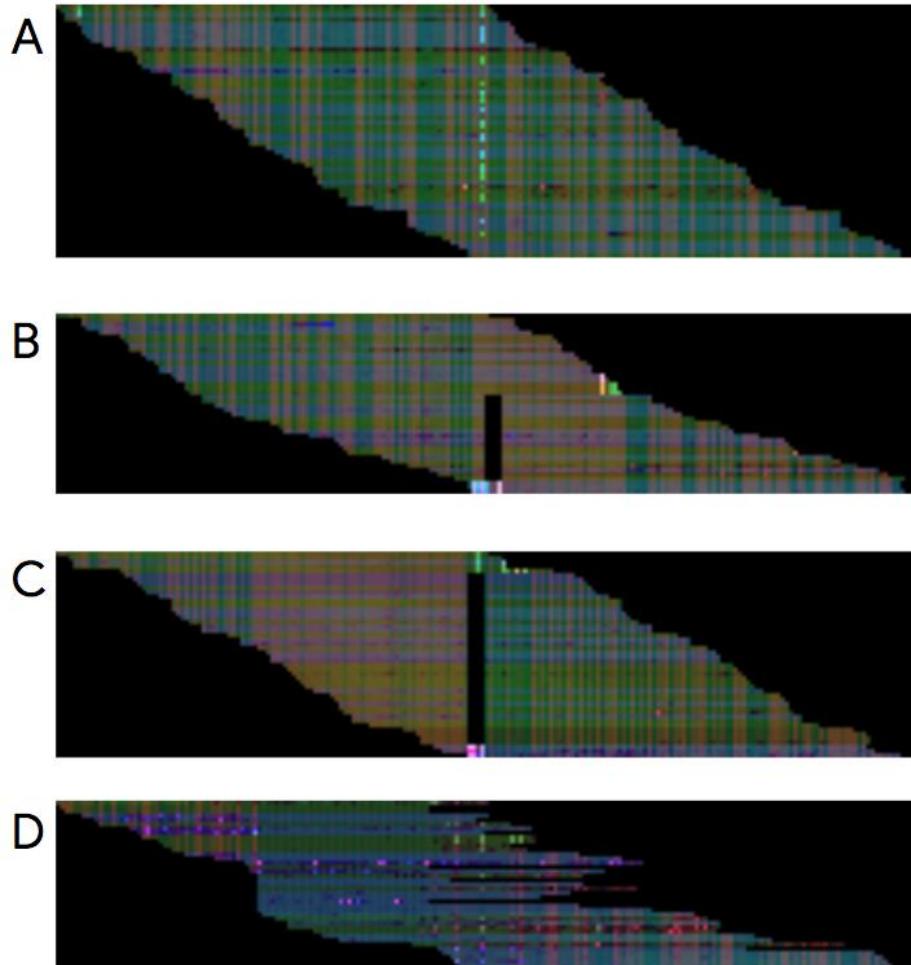
Model-based approach hinders detection of unknown types of complex SV

Challenges of detecting complex SV with long reads (2)



Repetitive regions complicate detection of SV

❖ Deep learning in variant detection



DeepVariant

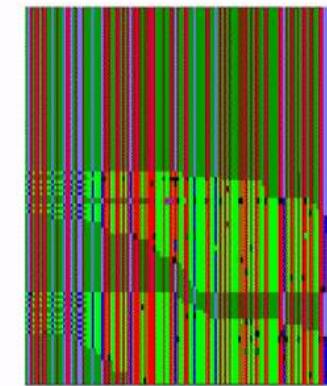
True SNP on one chromosome pair.

Deletion on one chromosome.

Deletion on both chromosomes.

False variant caused by errors.

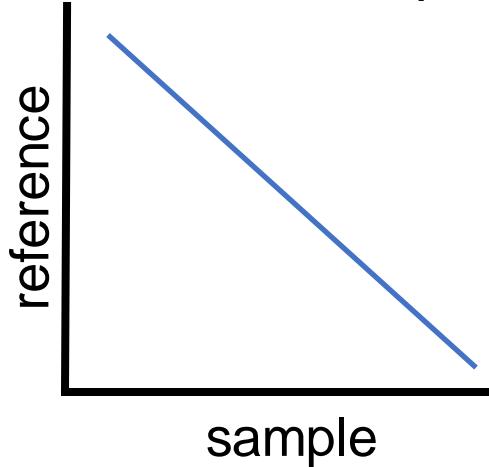
Read A
Read C
Read G
Read T
DELETION
Reference A
Reference C
Reference G
Reference T
INSERTION
Read Reverse Strand
Mate Reverse Strand
First in Pair
Second in Pair
Fails QC
PCR or Optical Duplicate
Mapping Quality



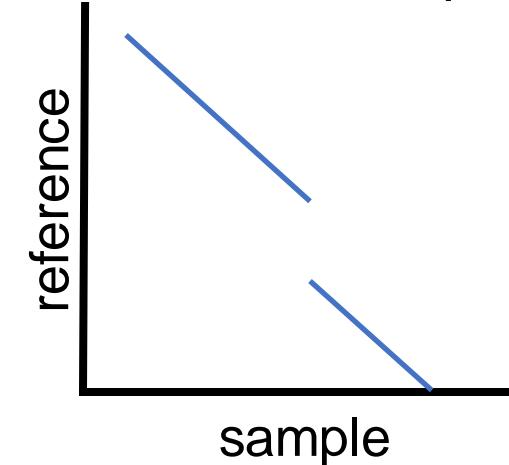
The sequence-to-image coding strategy of DeepVariant does not work on SV

Coding sequence similarity

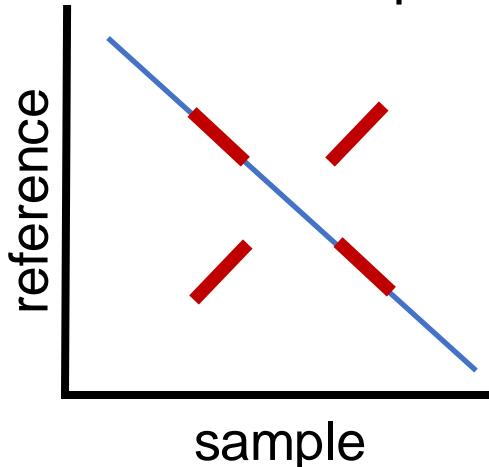
Identical & no repeat



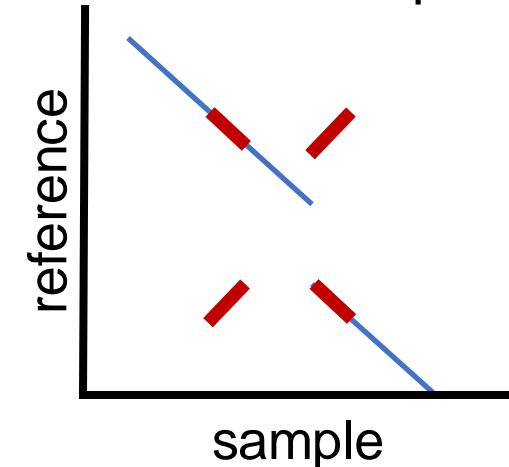
Different & no repeat



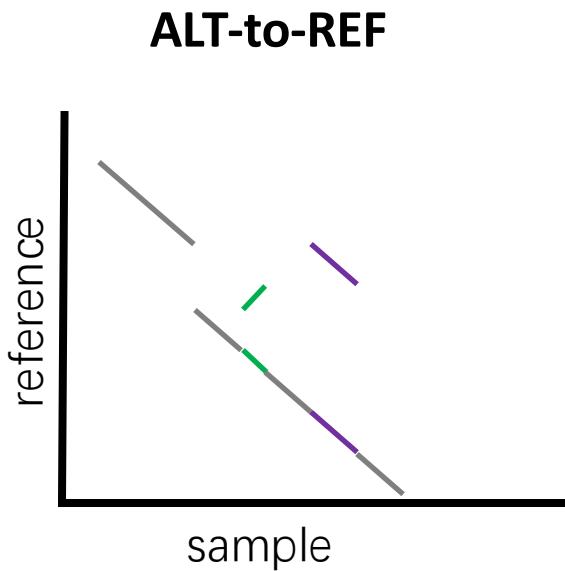
Identical & repeat



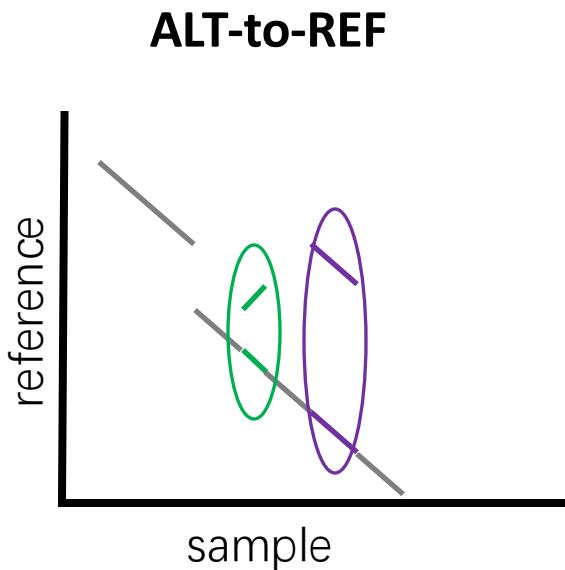
Different & repeat



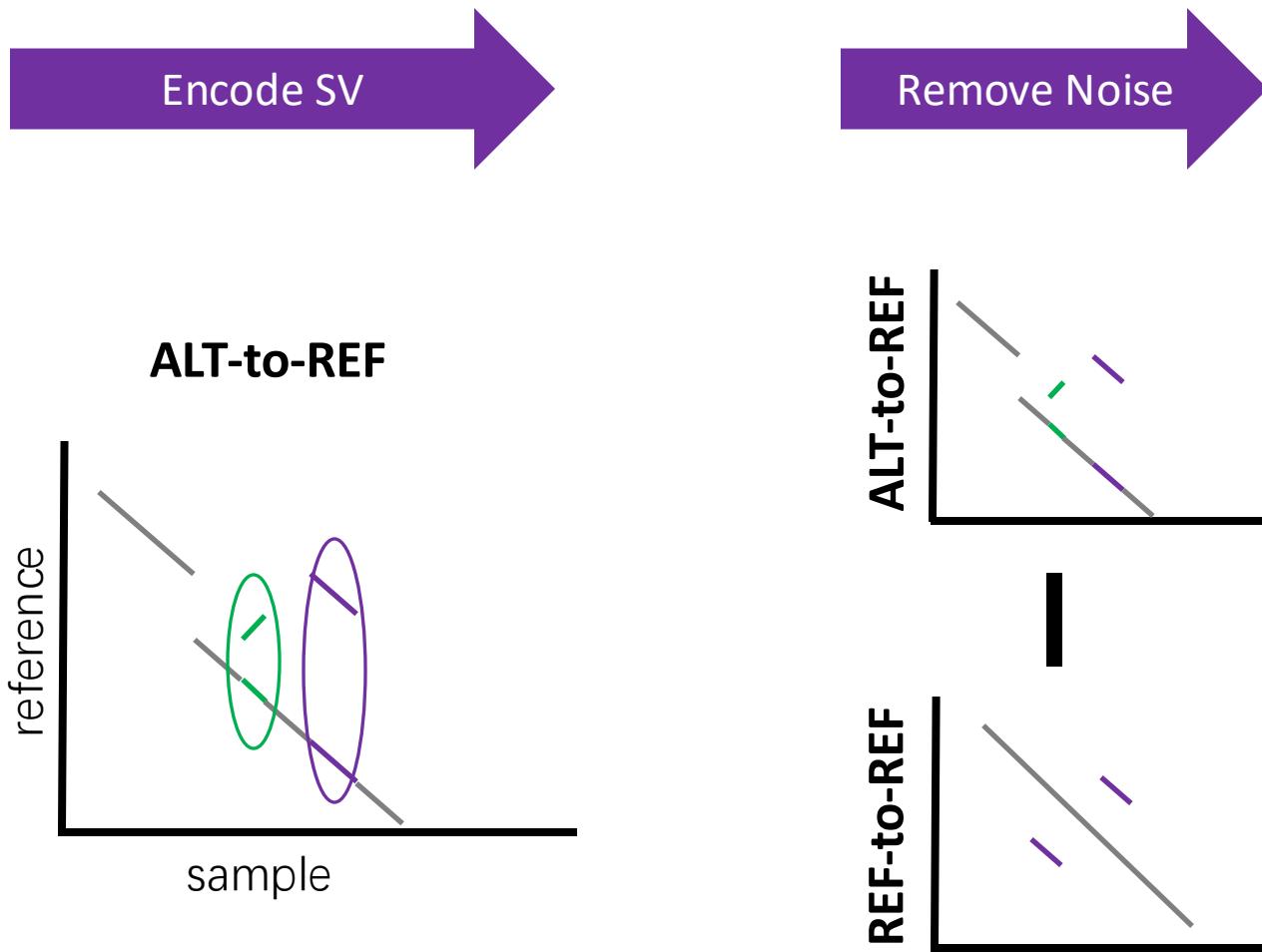
Coding sequence similarity with denoised dotplot



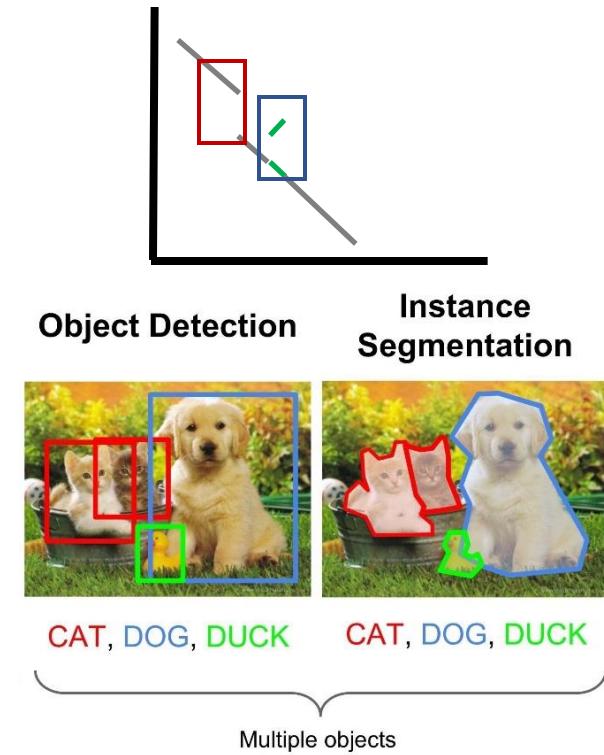
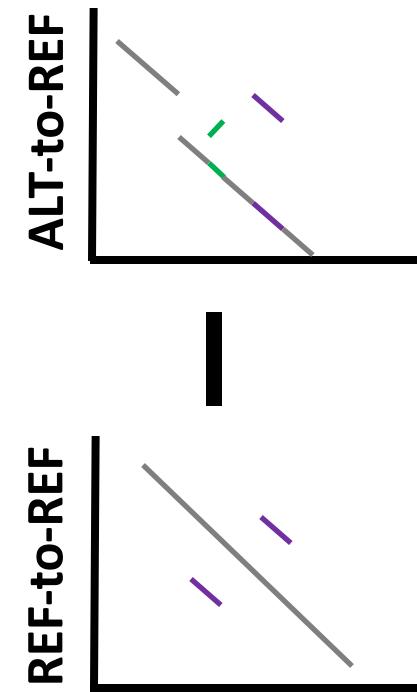
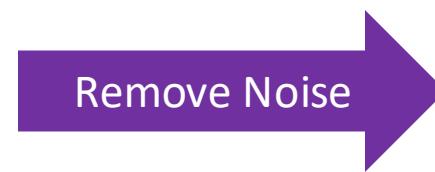
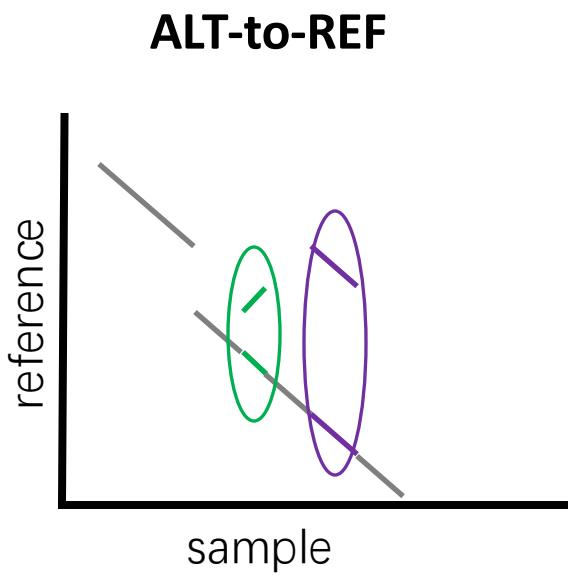
Coding sequence similarity with denoised dotplot



Coding sequence similarity with denoised dotplot



Coding sequence similarity with denoised dotplot



SVision workflow

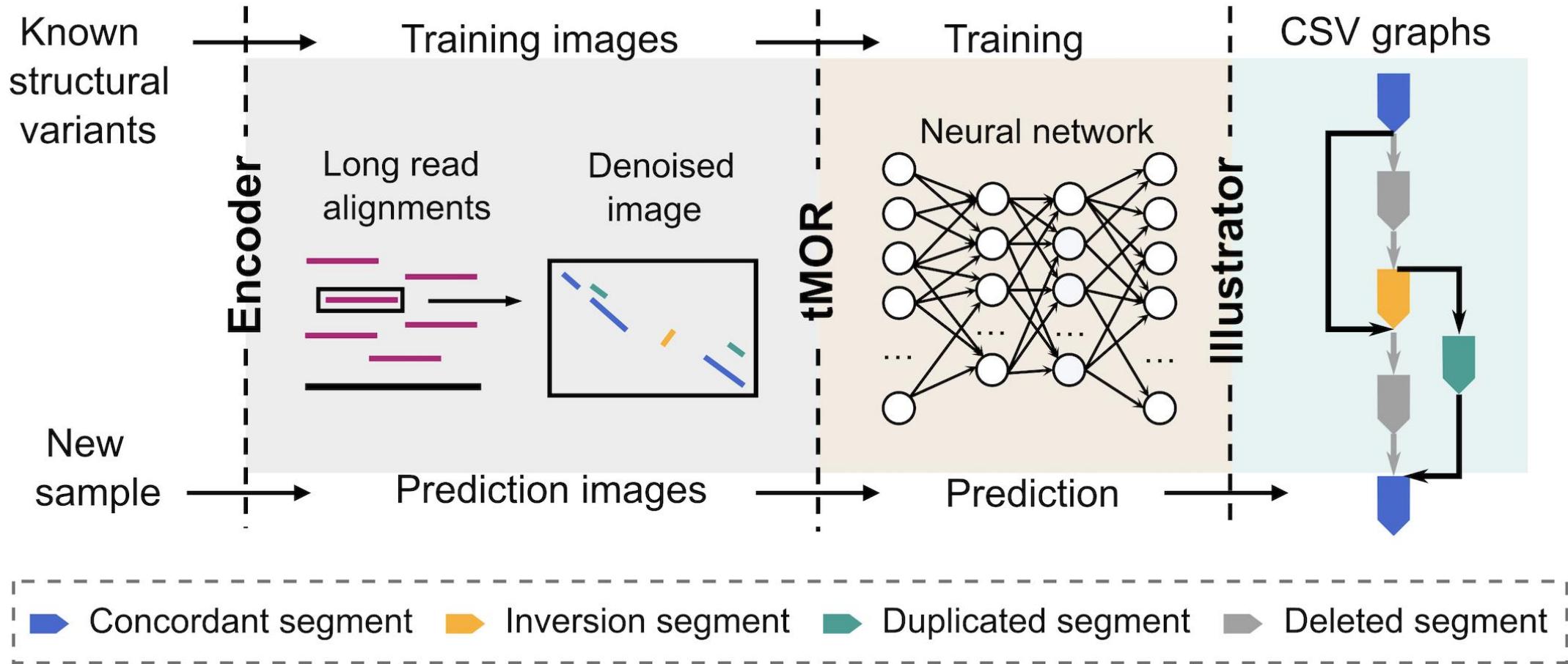


Image representation and object detection

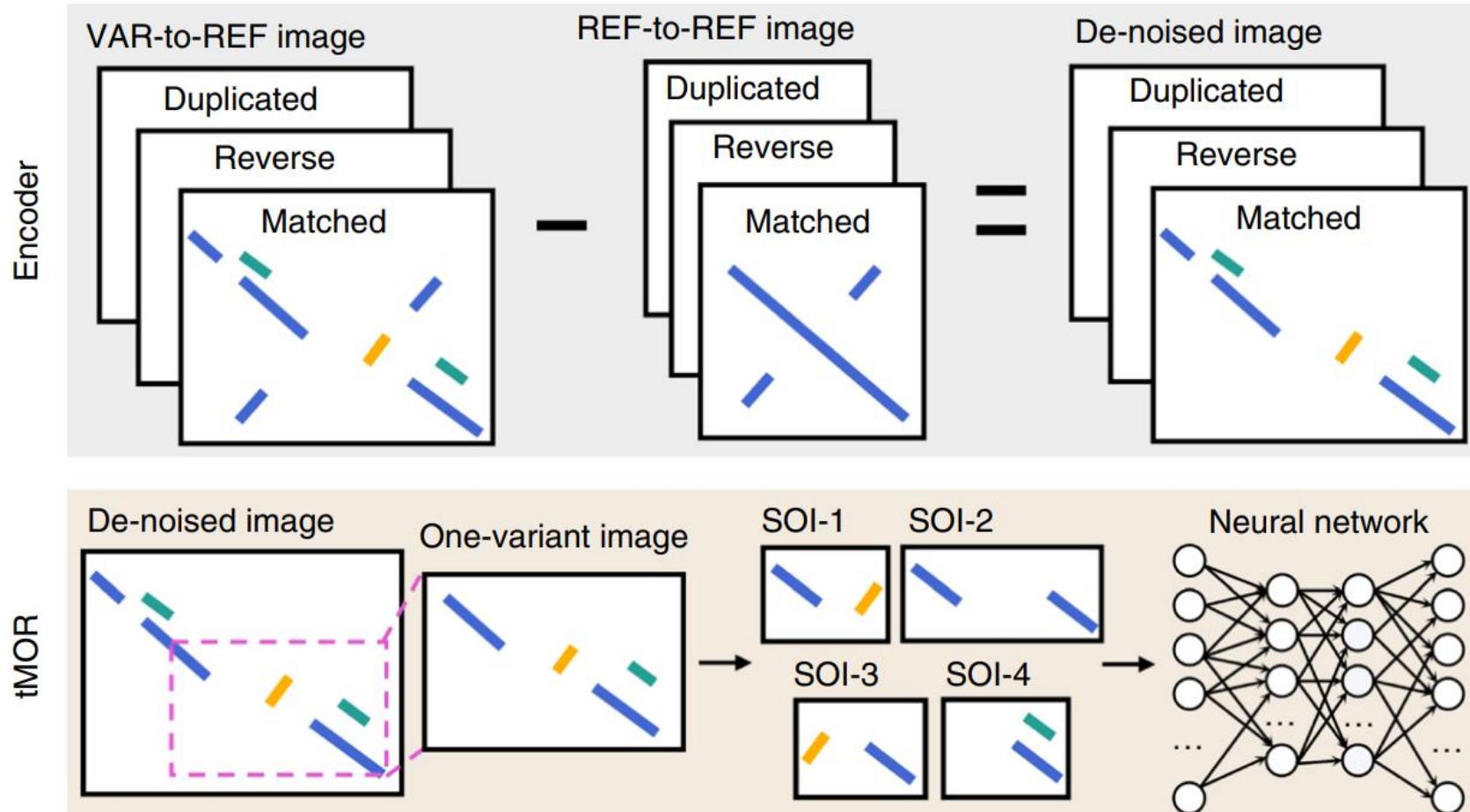
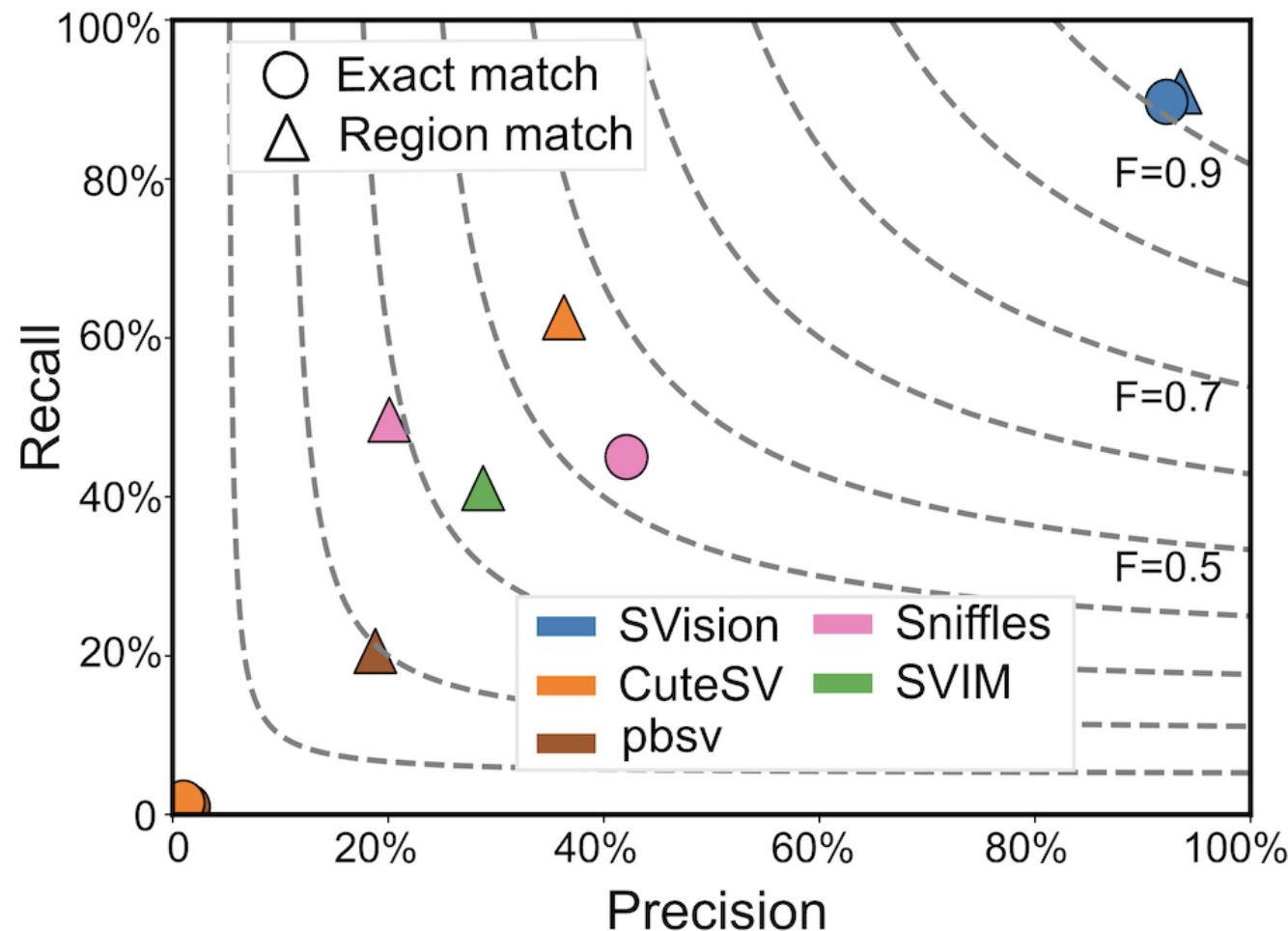


Image representation
and de-noising

Object detection

Five basic structural variant types

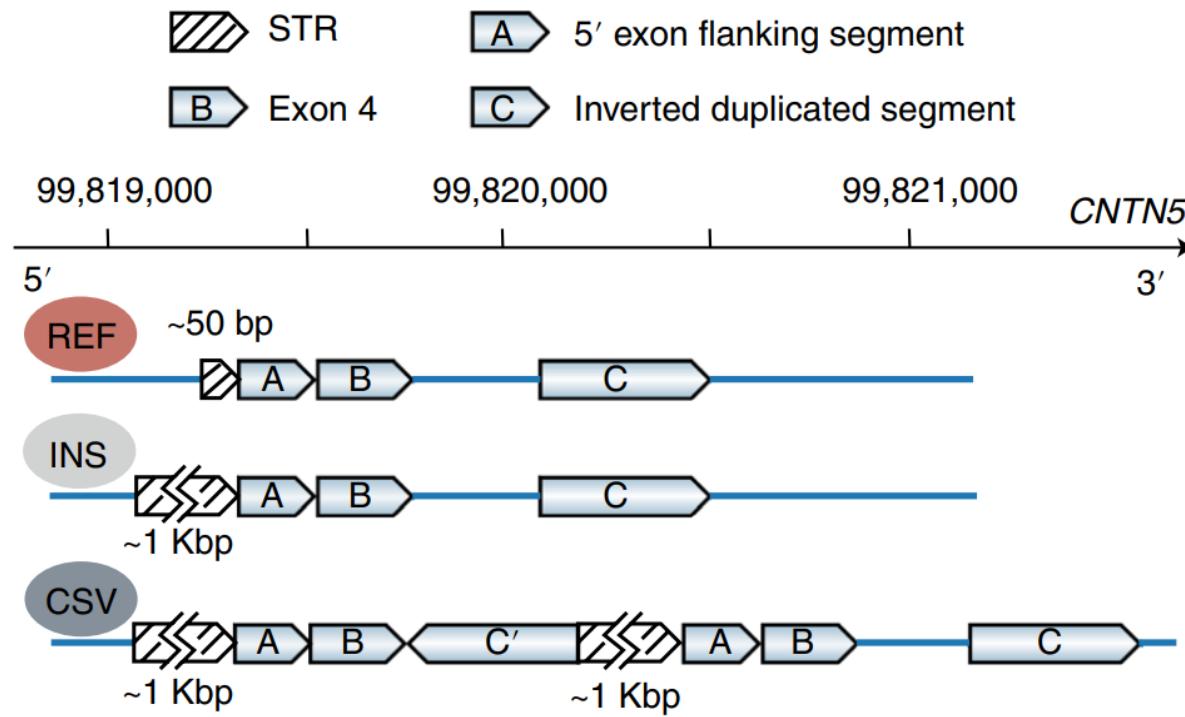
Calling simulated complex SV



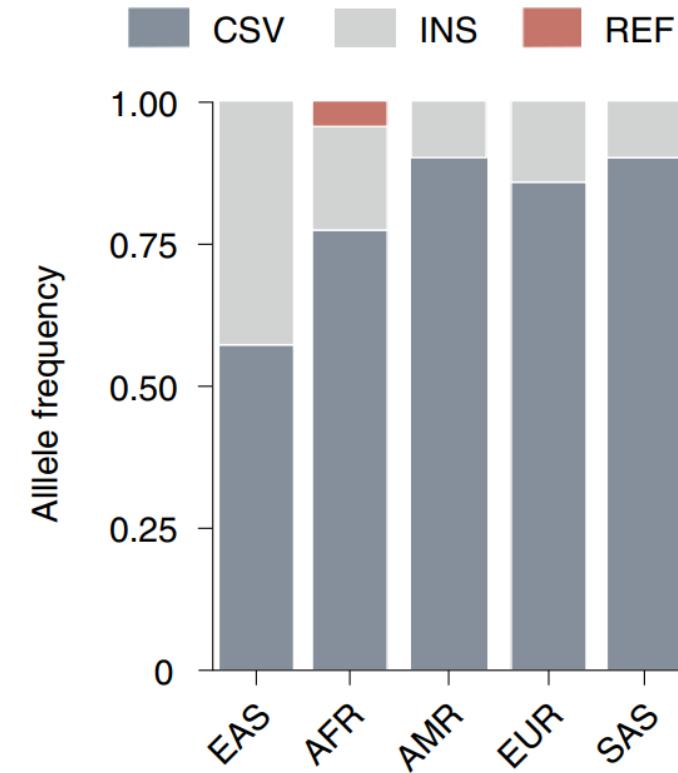
- Region match: correct discovery a CSV site;
- Exact match: correct discovery of CSV site and its subcomponents

❖ Resolving gene-related complex SV

Resolved CSV in CNTN5 coding region

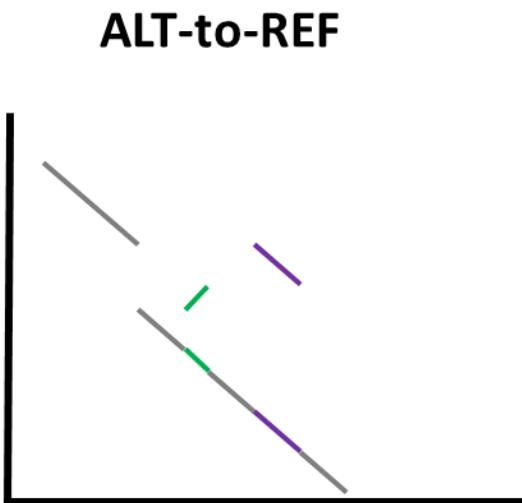


Population diversity

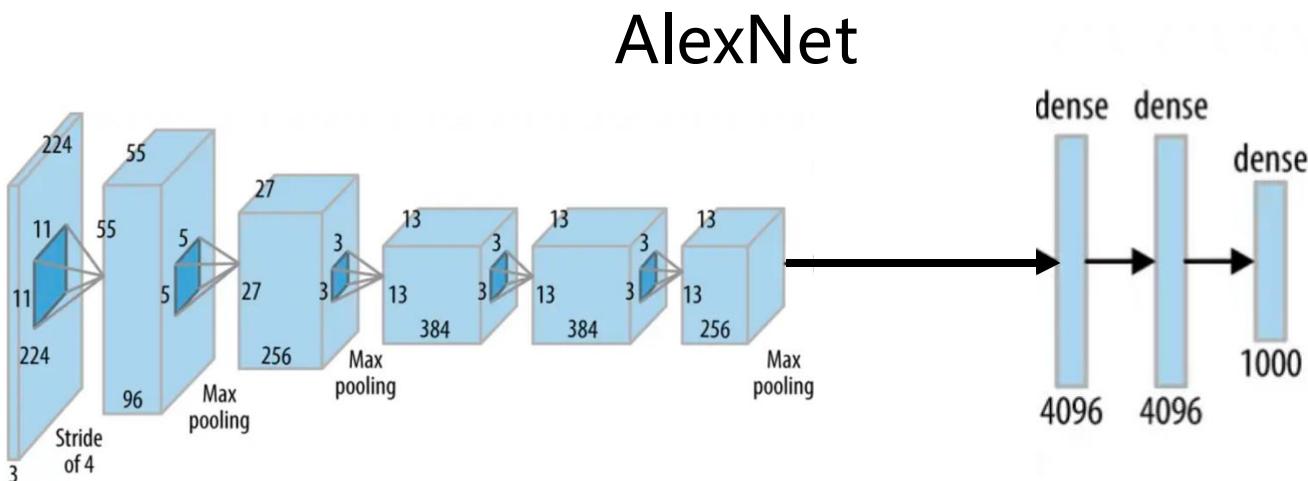


❖ Where we did well and how to improve?

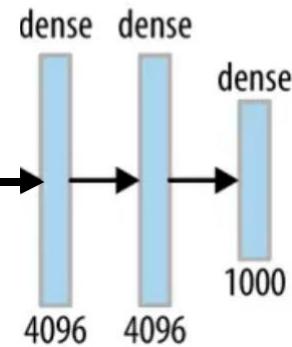
Sequence to
image



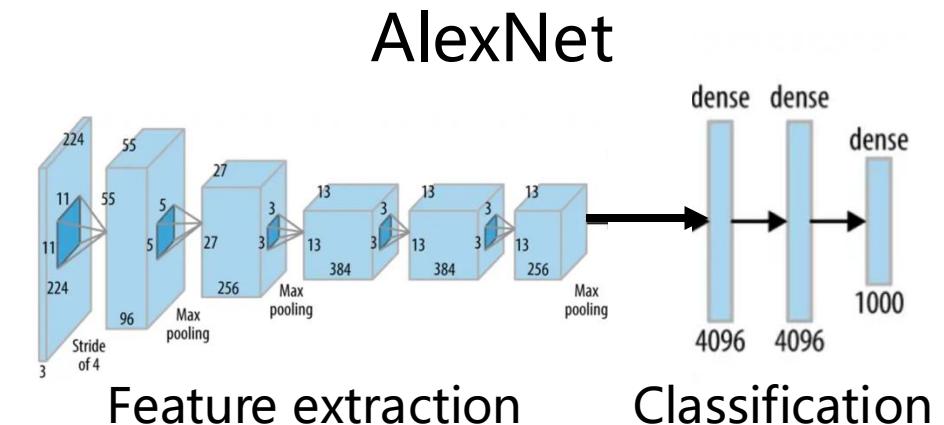
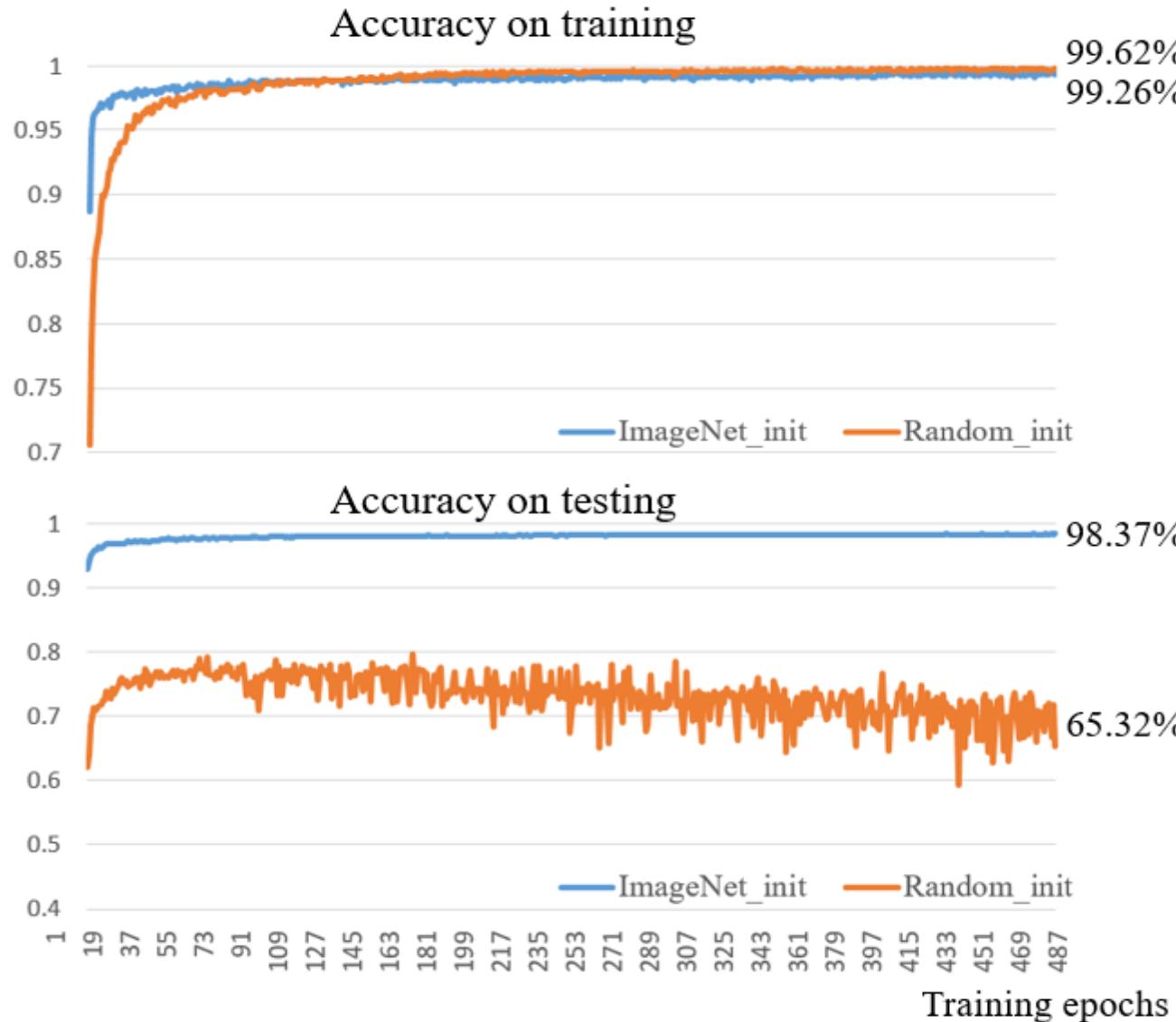
Feature
extraction



Classification



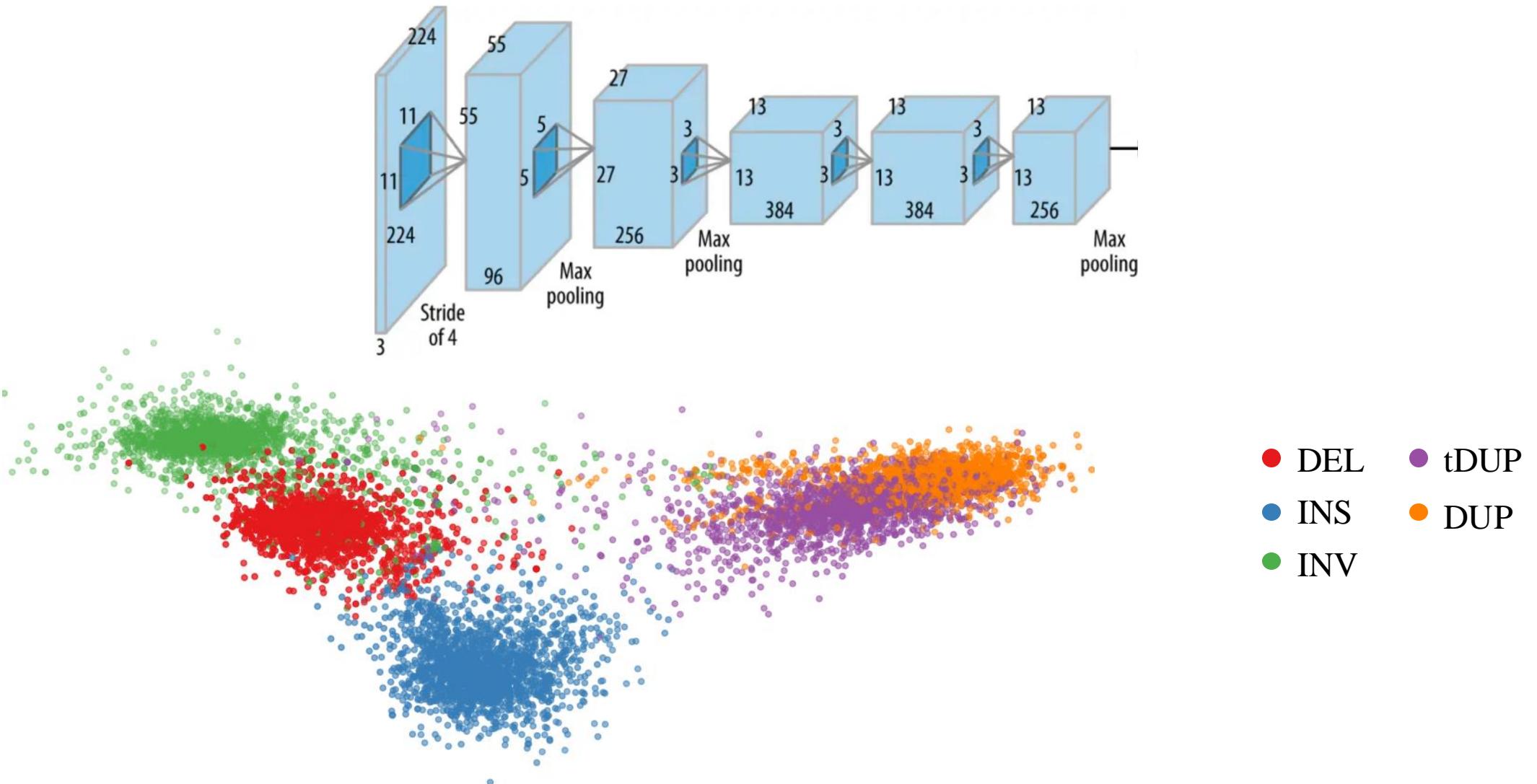
❖ Training AlexNet for SVision



- **ImageNet_init** = set 5 feature extraction layers as ImageNet setting, only train 3 classification layers
- **Random_init** = train all 8 layers (5 feature extraction + 3 classification layers)

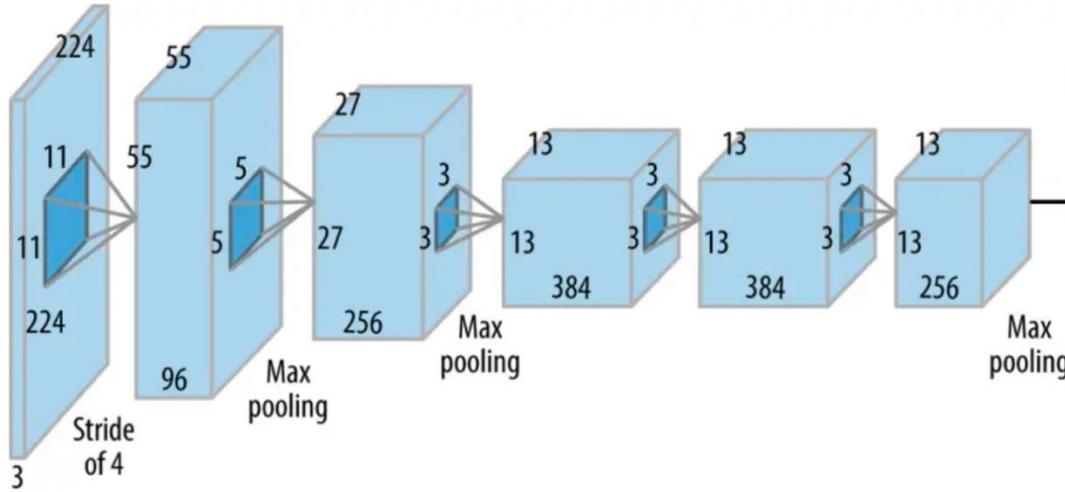
❖ AlexNet with ImageNet_init setting extracts fine features from training set

Feature extraction



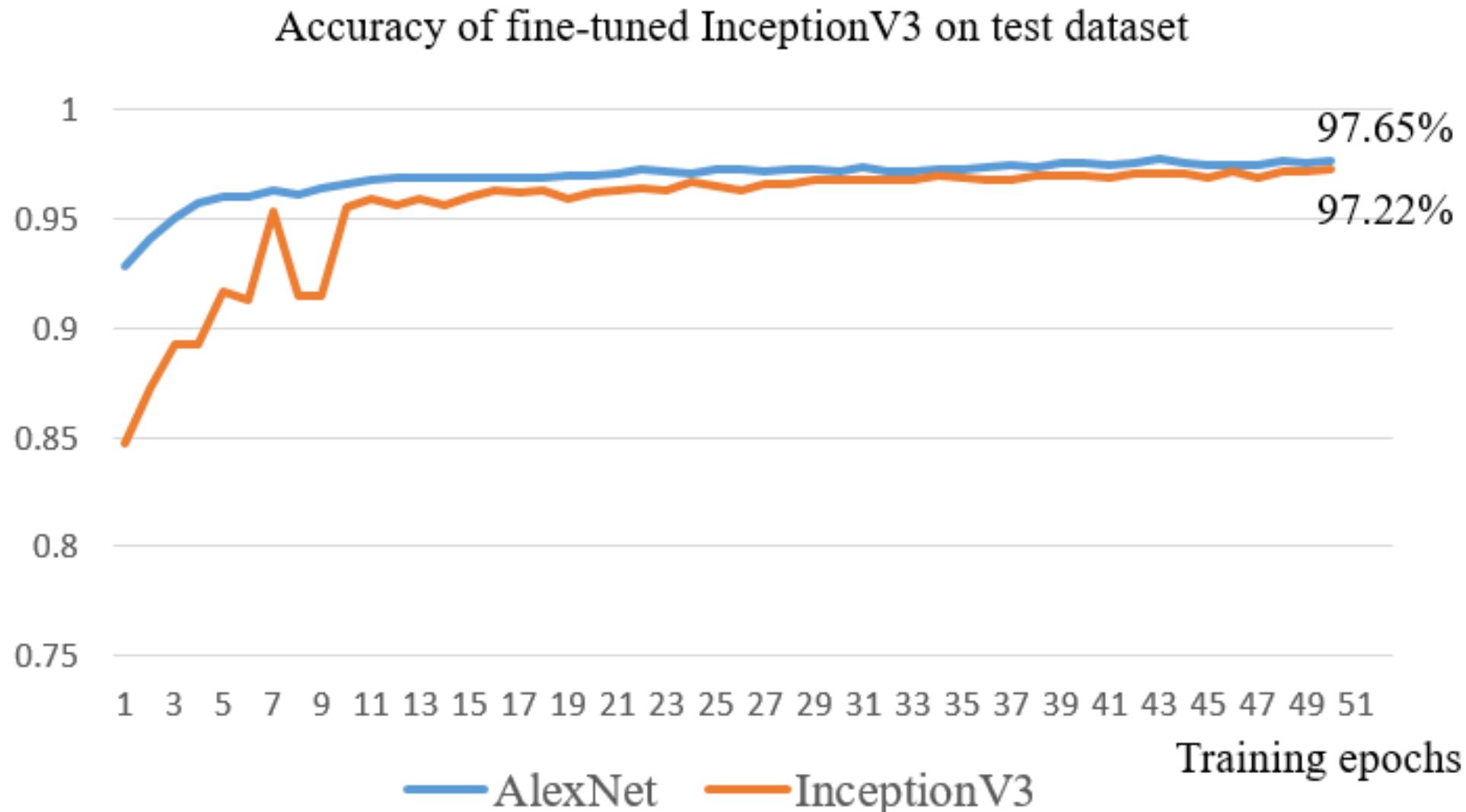
❖ Features extracted by AlexNet with ImageNet_init fit other ML approaches

Feature extraction

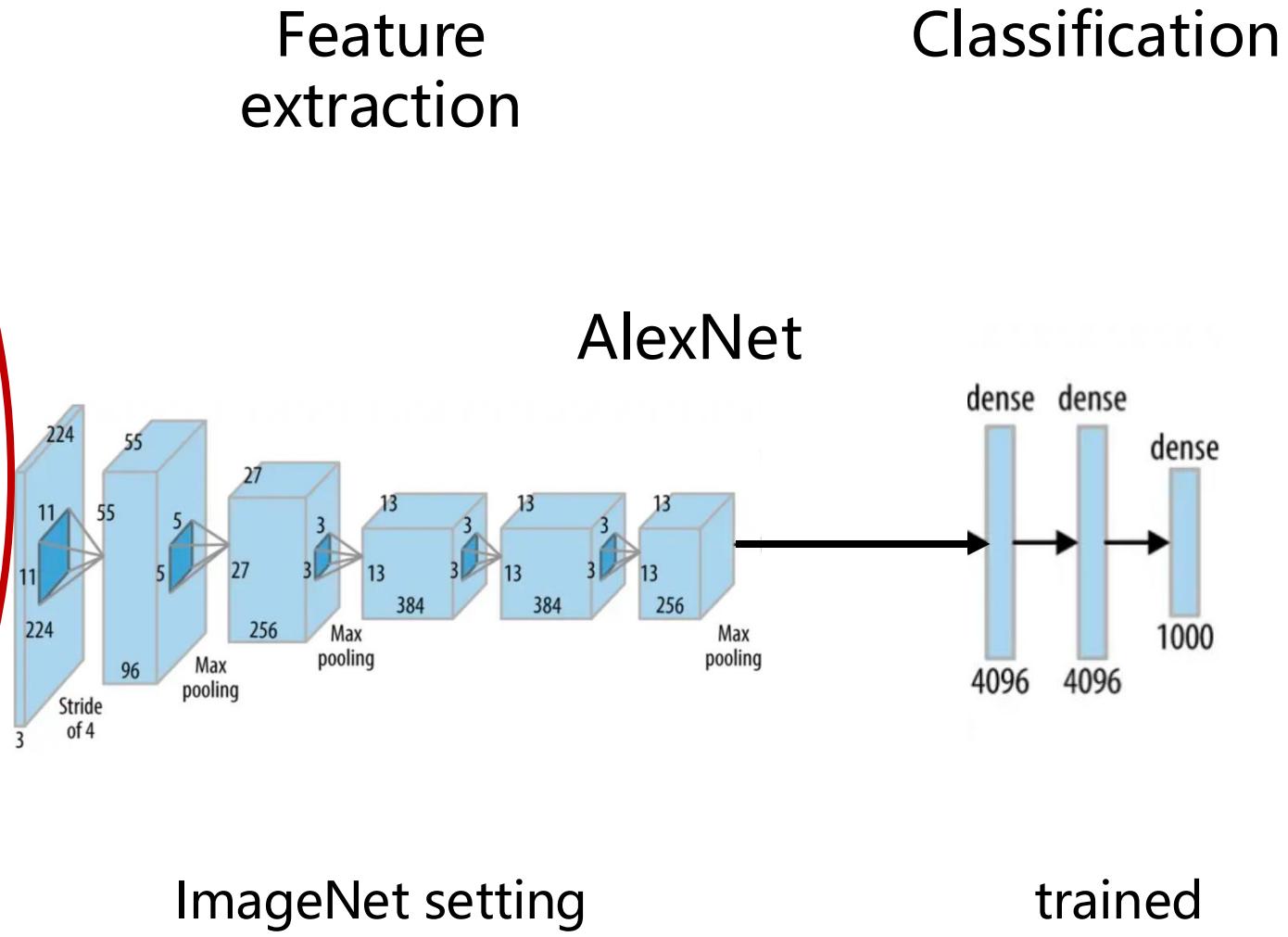
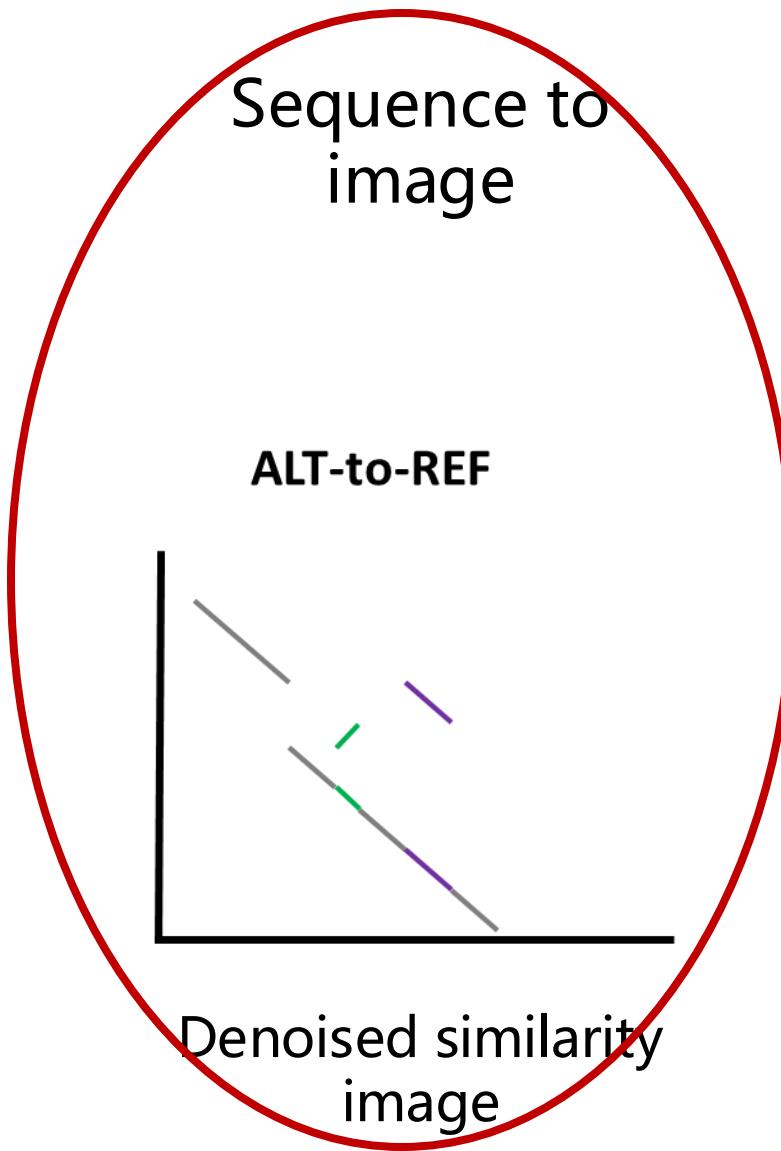


Methods	Recall	Precision
Logistic regression	88.49%	91.14%
SVM	97.12%	97.09%
Random forests	96.59%	96.77%

❖ Other CNN framework works too



❖ Where we did well and how to improve?





Check for updates

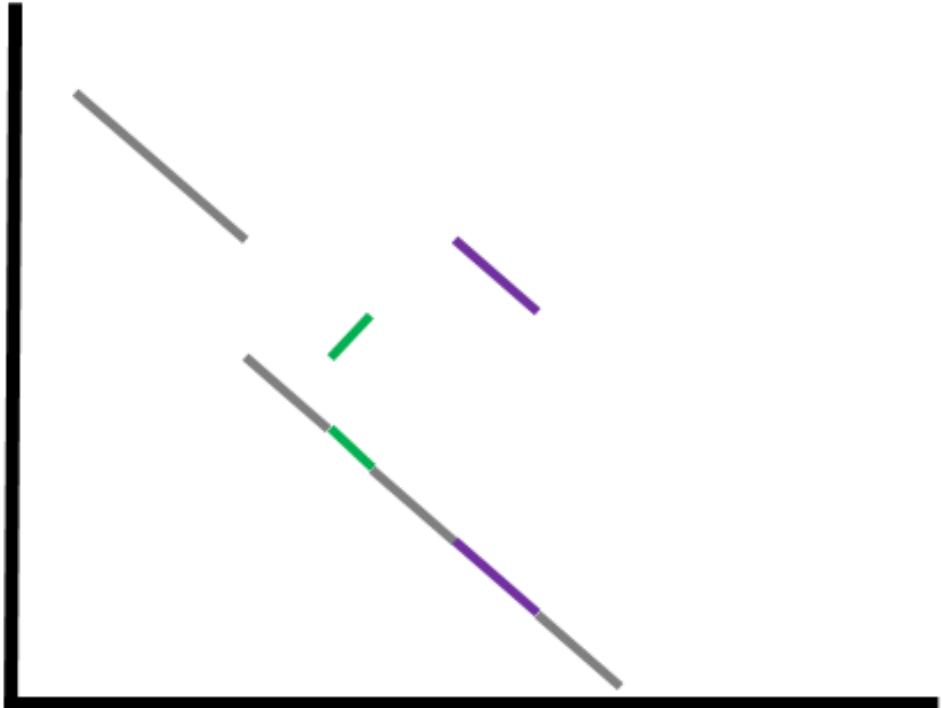
SVision: a deep learning approach to resolve complex structural variants

Jiadong Lin  ^{1,2,3,4,13}, Songbo Wang  ^{1,2,3,13}, Peter A. Audano  ⁵, Deyu Meng ^{1,6,7}, Jacob I. Flores ⁵, Walter Kosters  ⁴, Xiaofei Yang  ^{1,8}, Peng Jia ^{1,2,3}, Tobias Marschall  ⁹, Christine R. Beck  ^{5,10} and Kai Ye  ^{1,2,3,11,12} 

Where we did well and how to improve?

SVision Seq2image

Compare ref and sample



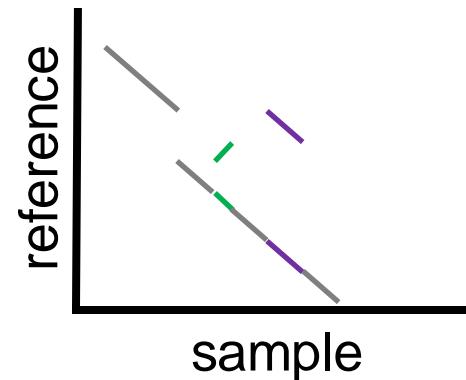
SVision Seq2image encodes

- Mapping information
- Sequence features
- Domain knowledge

New application scenarios: SVision-pro

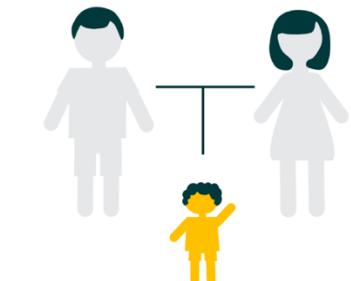
Single genome variant calling

Sample vs Reference
germline



Multi-genome comparison

Sample vs Sample
De novo or somatic



De novo variants

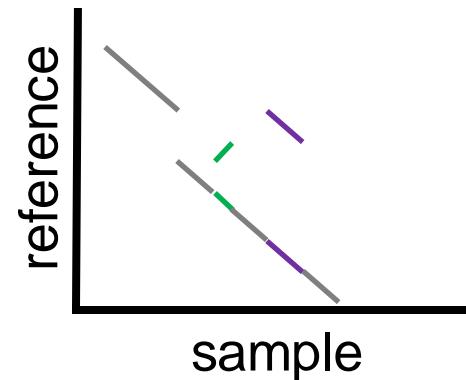


Somatic variants

New application scenarios: SVision-pro

Single genome variant calling

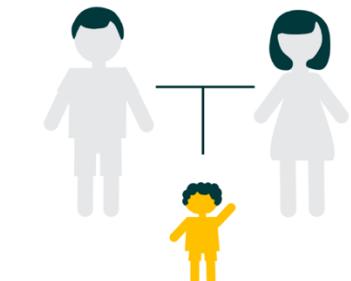
Sample vs Reference
germline



$\text{Diff}(\text{reference}, \text{sample})$

Multi-genome comparison

Sample vs Sample
De novo or somatic



De novo variants



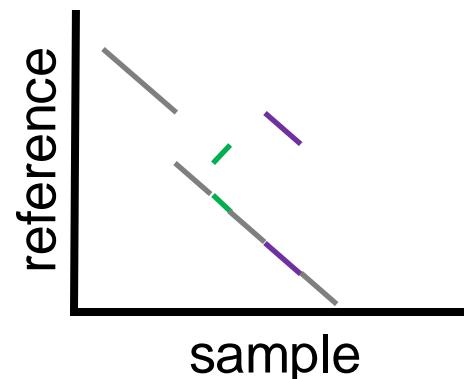
Somatic variants

$\text{Diff}(\text{Diff}(\text{reference}, \text{normal}), \text{Diff}(\text{reference}, \text{tumor}))$

New application scenarios: SVision-pro

Single genome variant calling

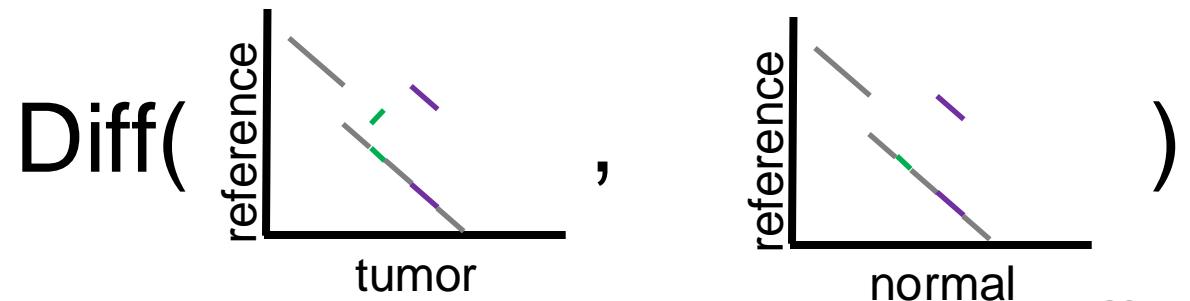
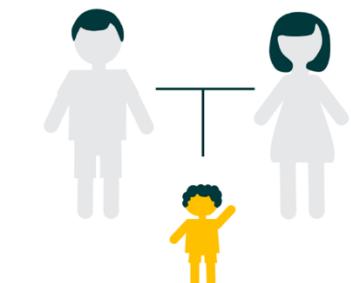
Sample vs Reference
germline



Diff(reference, sample)

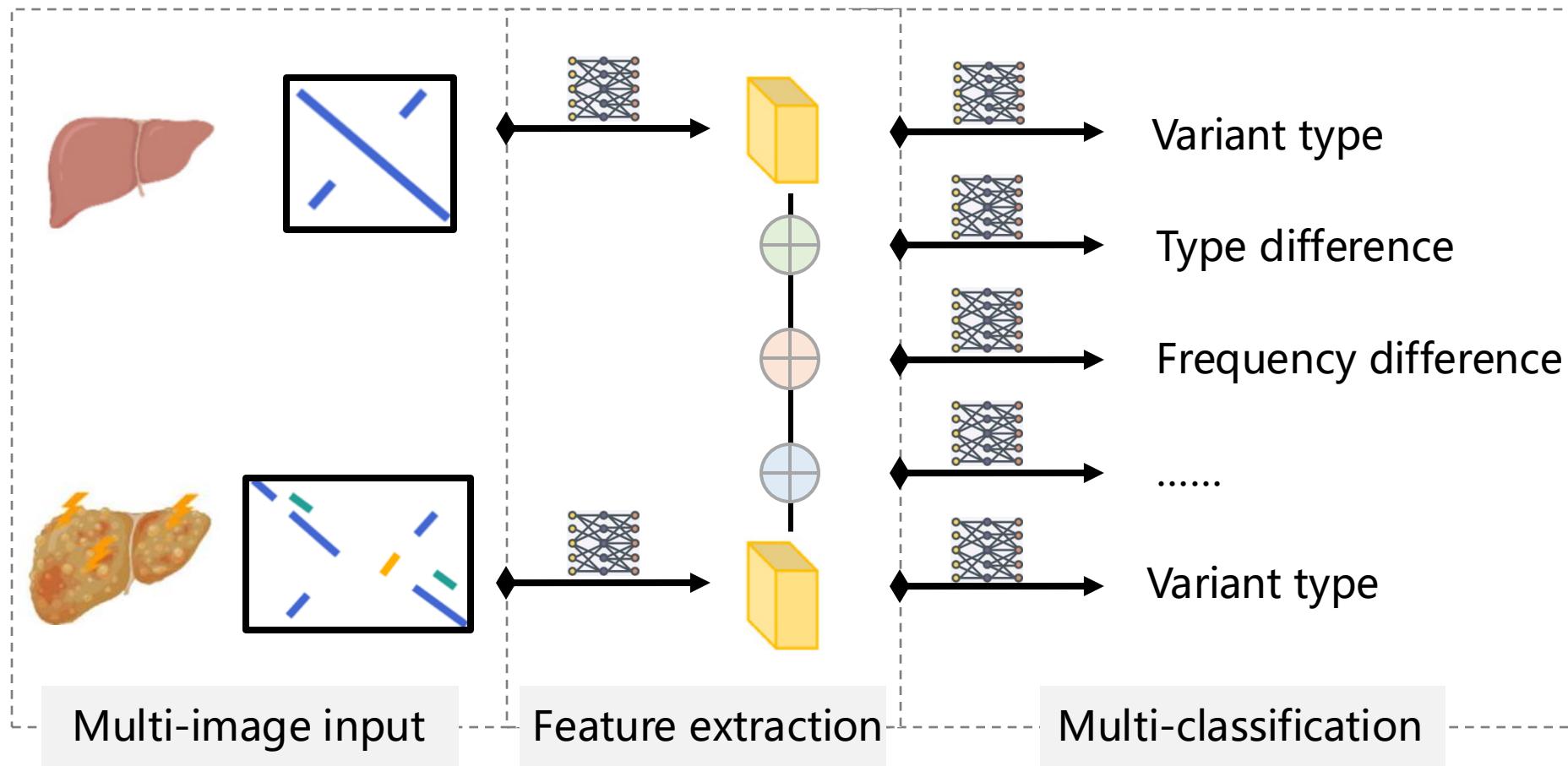
Multi-genome comparison

Sample vs Sample
De novo or somatic



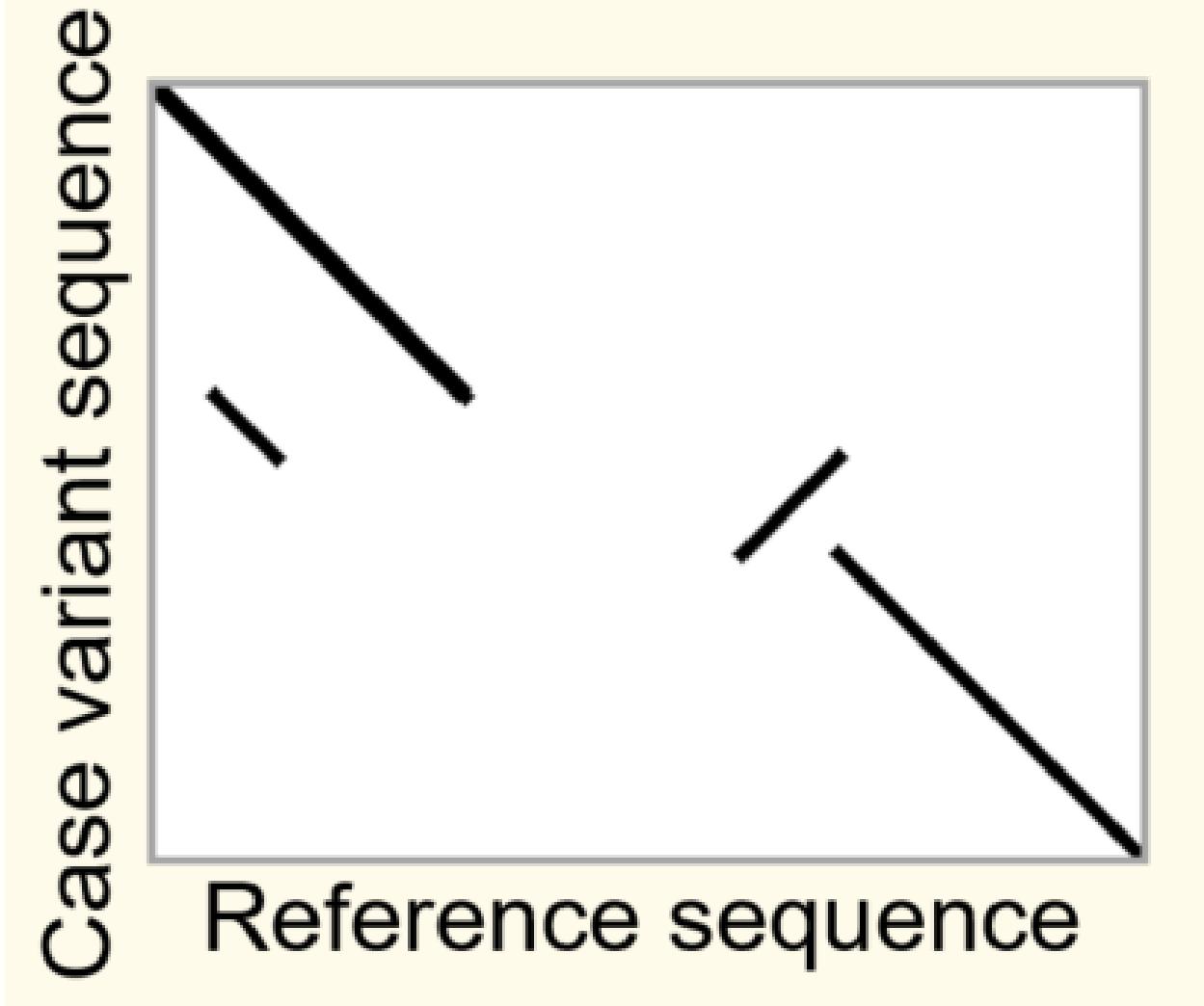


Do we directly compare two images? NO

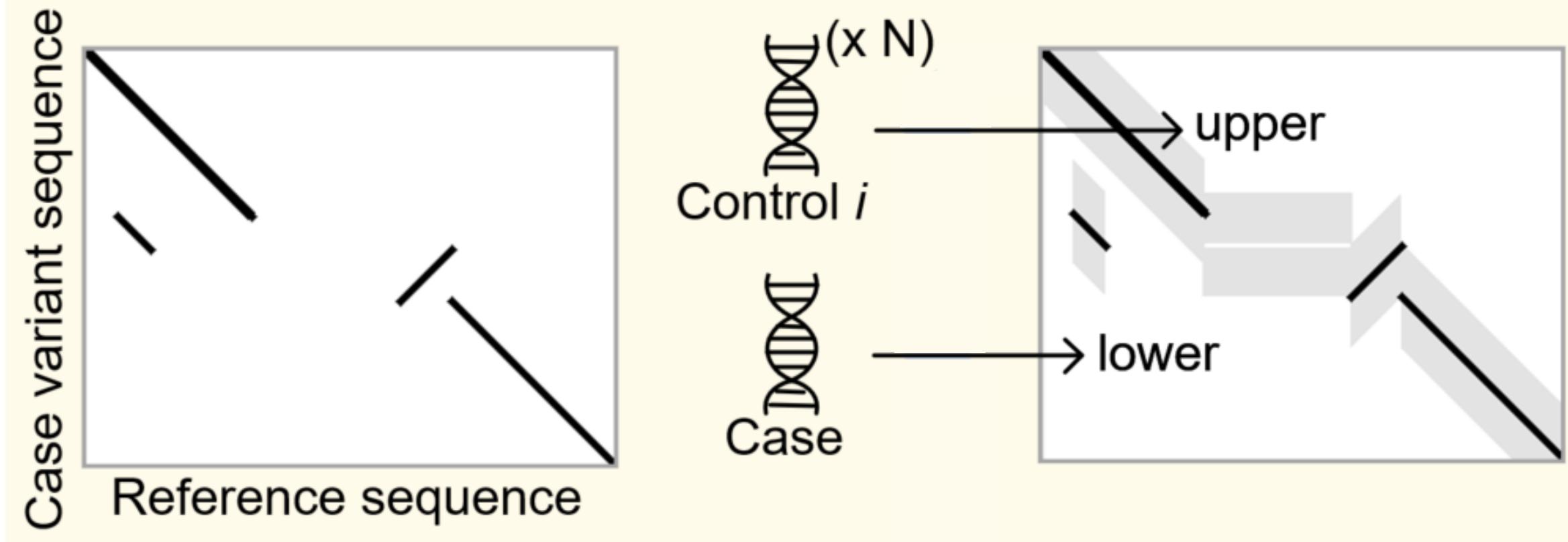


Abandon direct comparison due to computational burden
we need a lightweight and elegant approach

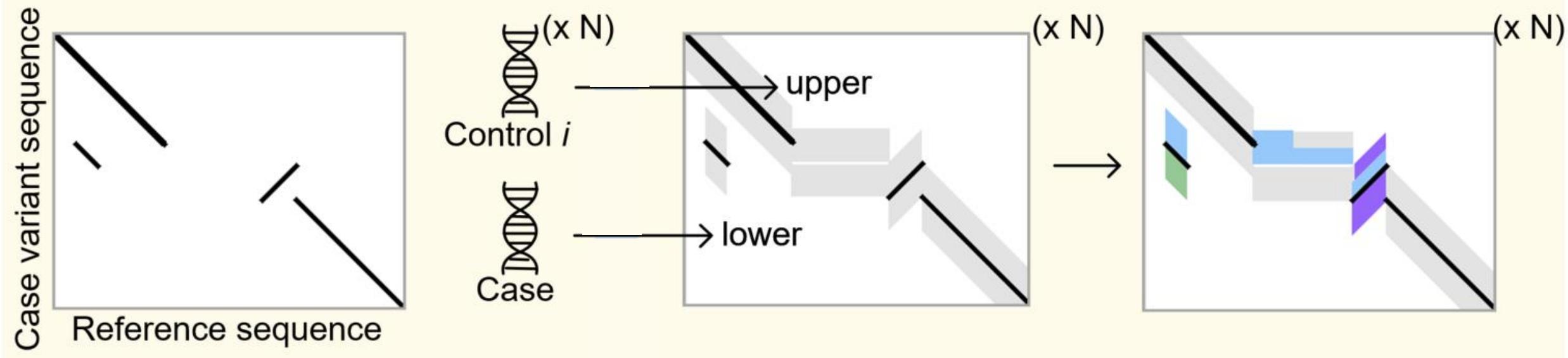
Our strategy: Genome-to-genome representation within one image



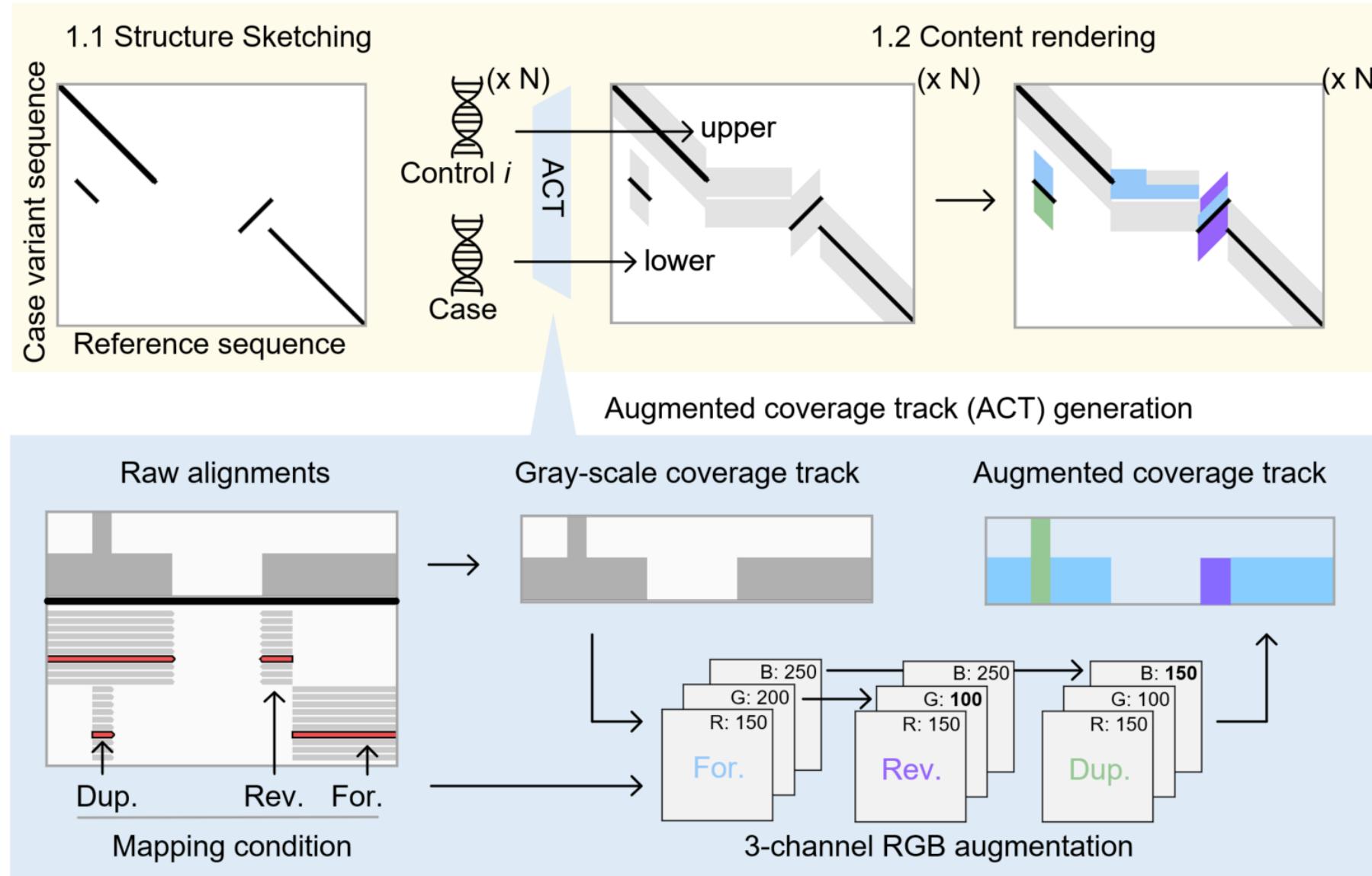
Our strategy: Genome-to-genome representation within one image



Our strategy: Genome-to-genome representation within one image



Our strategy: Genome-to-genome representation within one image



Flexible image setting

Flexible image properties for different sensitivity requirements



Image size: 256 x 256
Track height: 25
Min-representable AF: 0.04

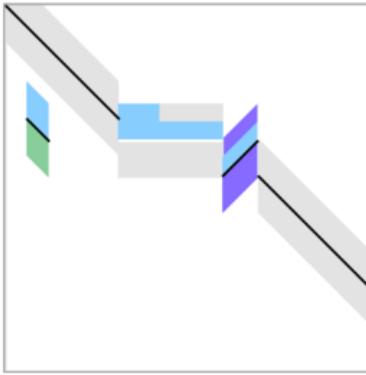


Image size: 512 x 512
Track height: 50
Min-representable AF: 0.02

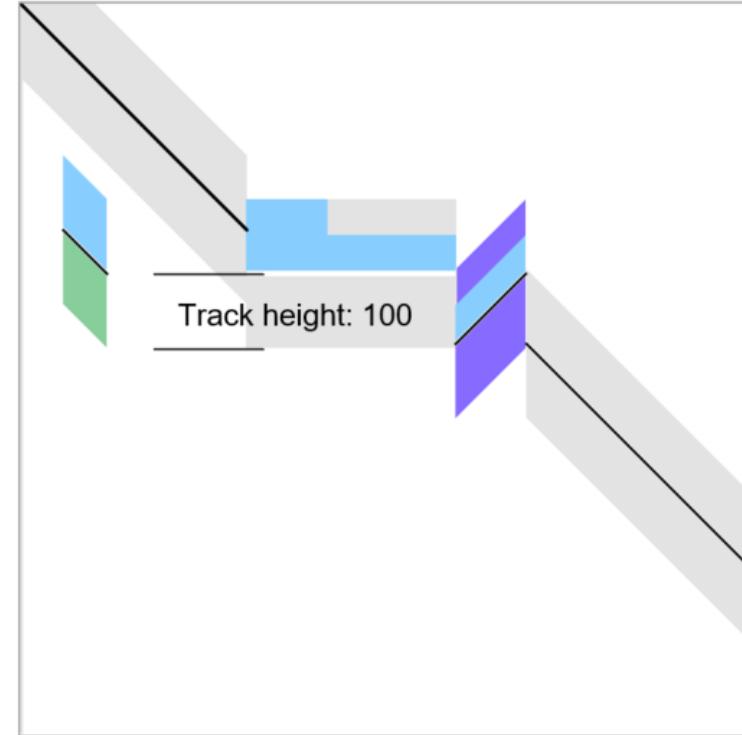


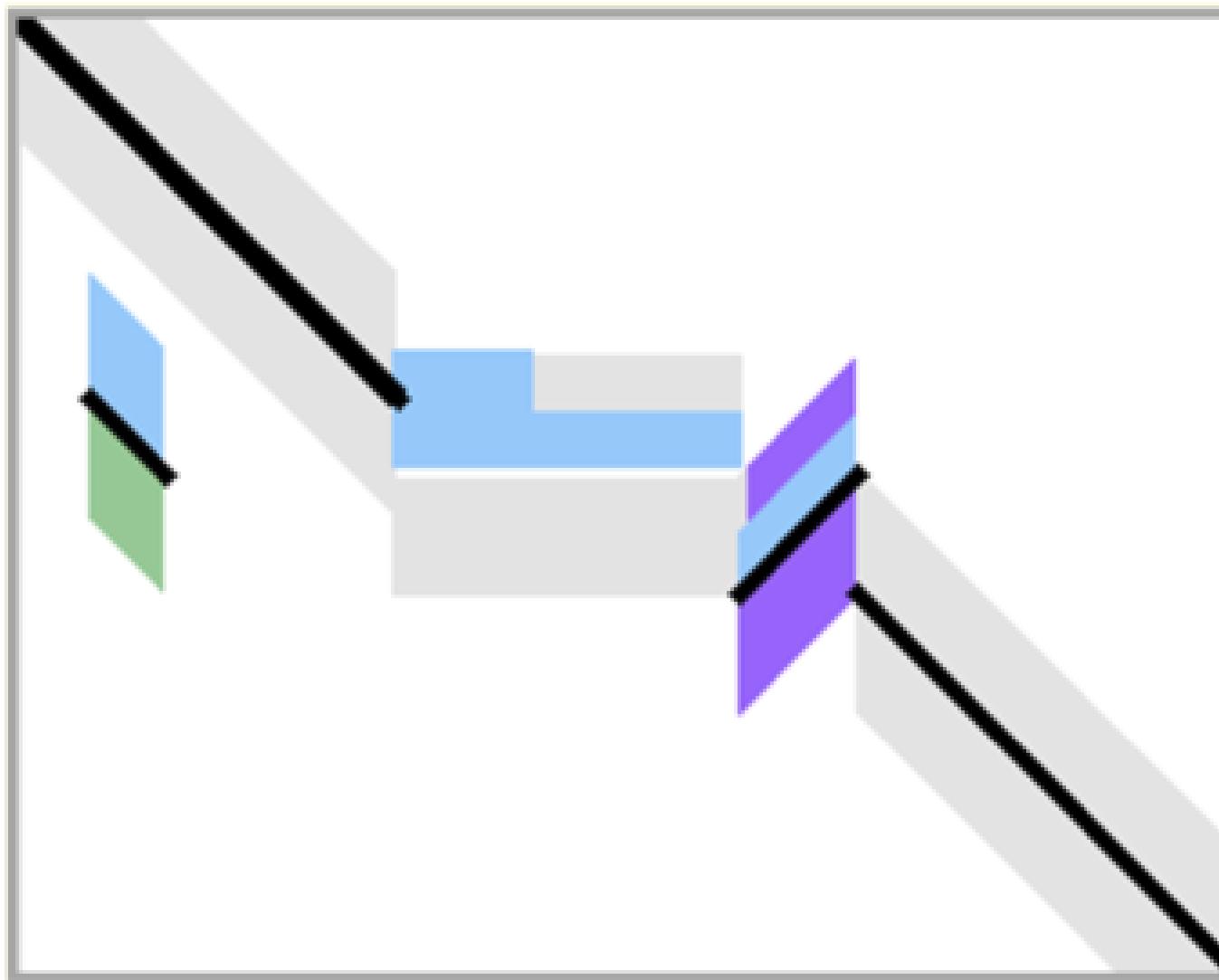
Image size: 1024 x 1024
Track height: 100
Min-representable AF: 0.01

.....
(customization)

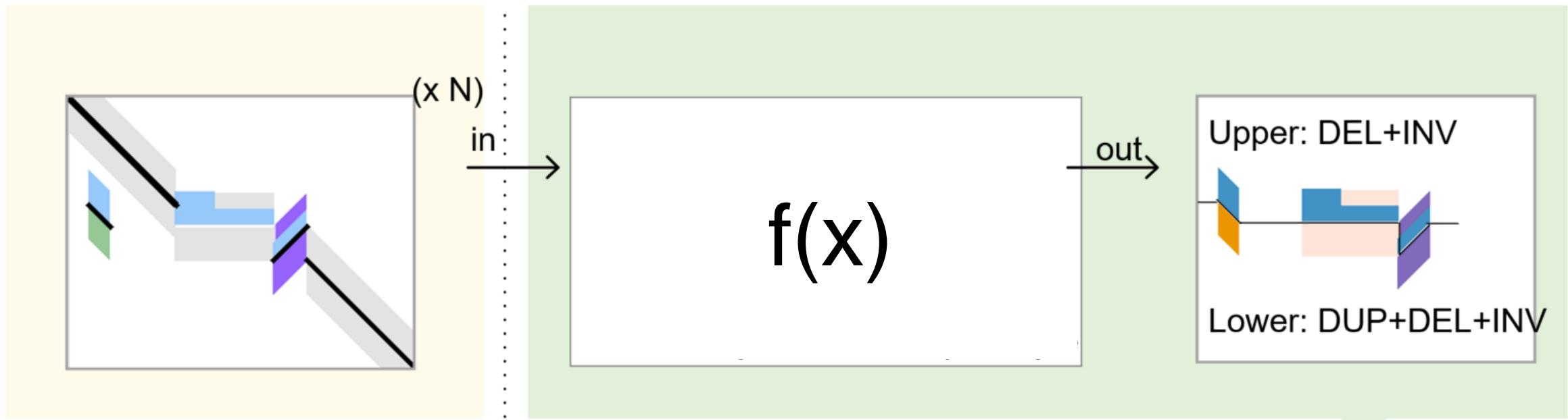
Image sizes
Track heights
Colors

Larger image size leads to lower AF

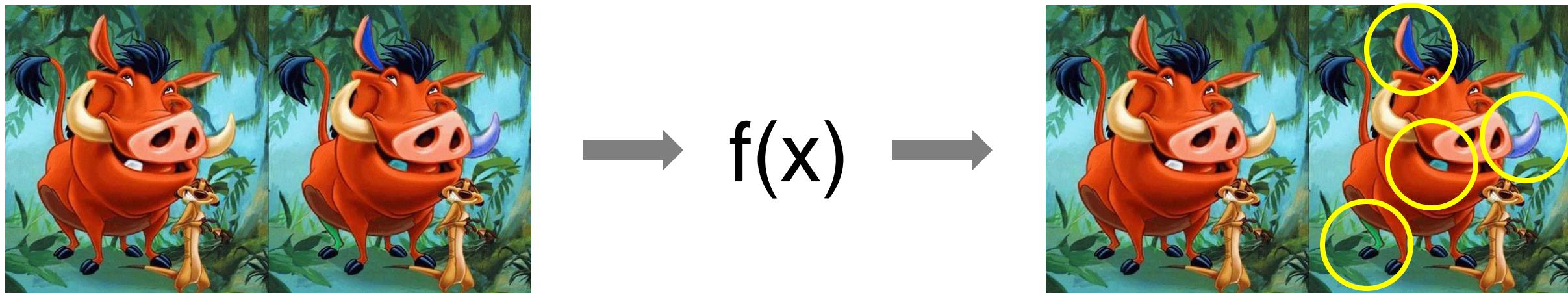
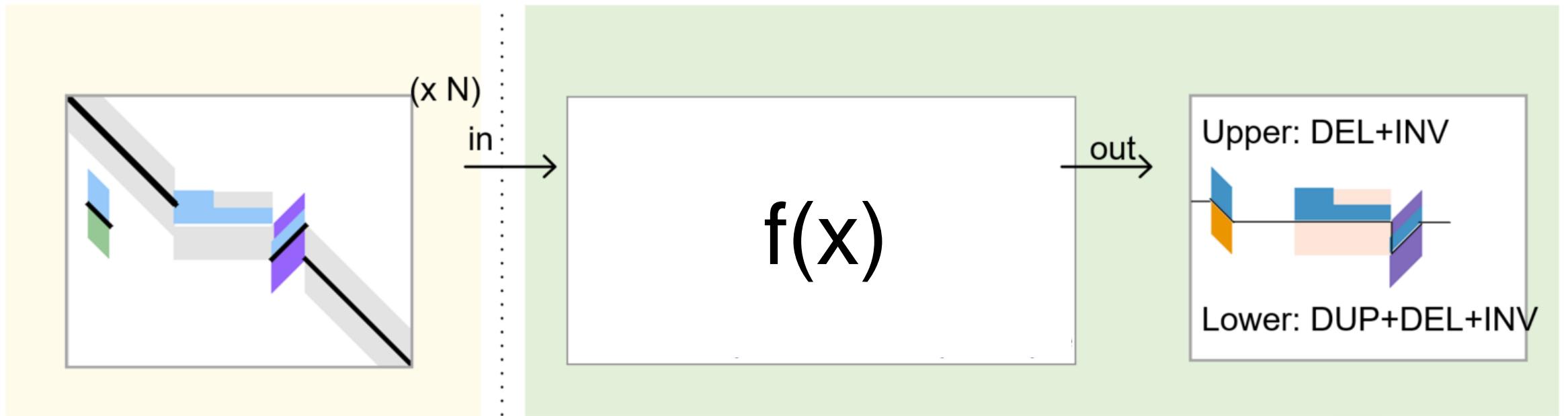
Comparative recognition module



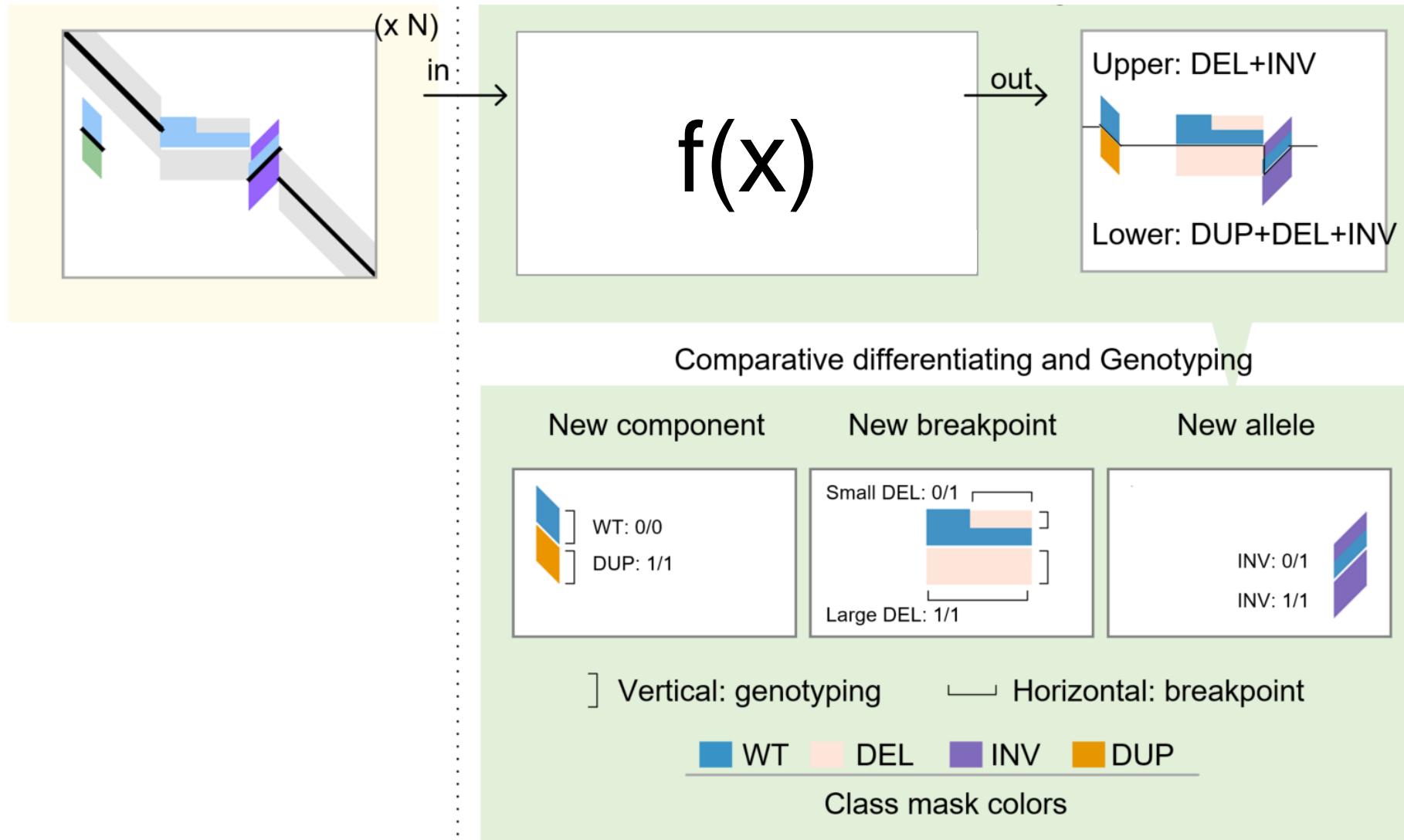
Comparative recognition module



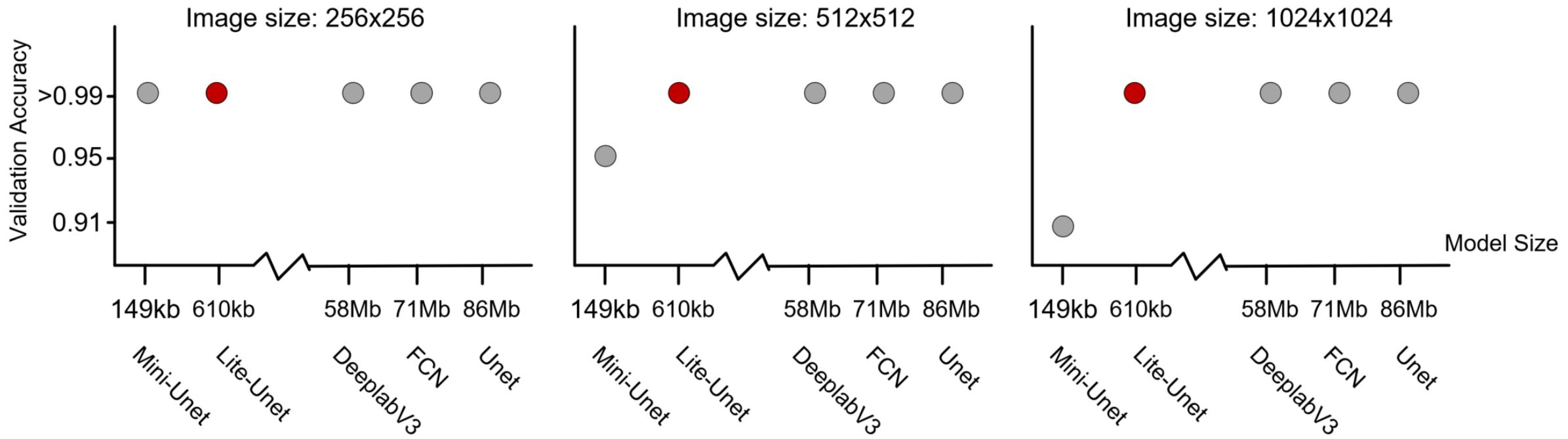
Comparative recognition module



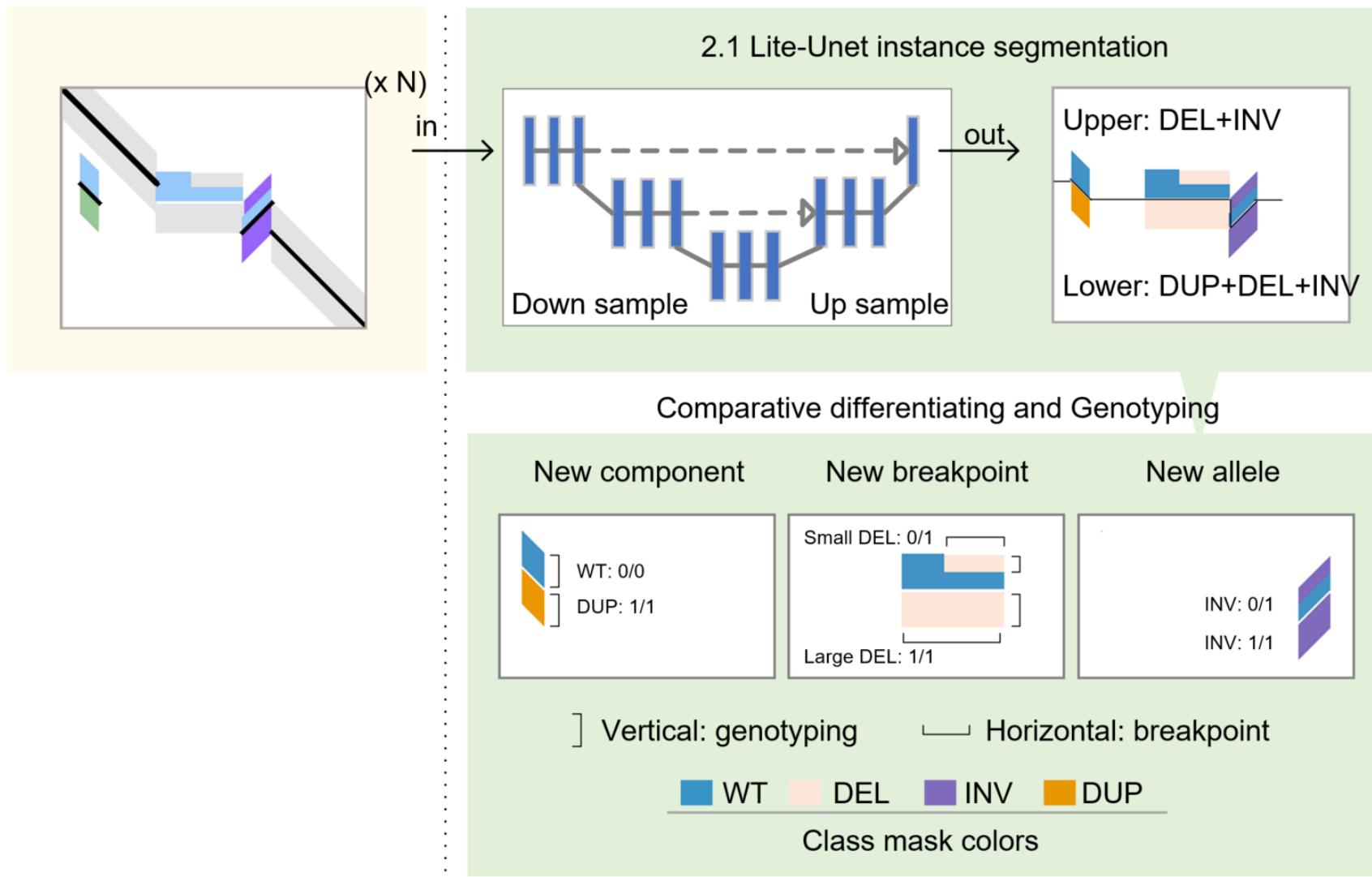
Comparative recognition module



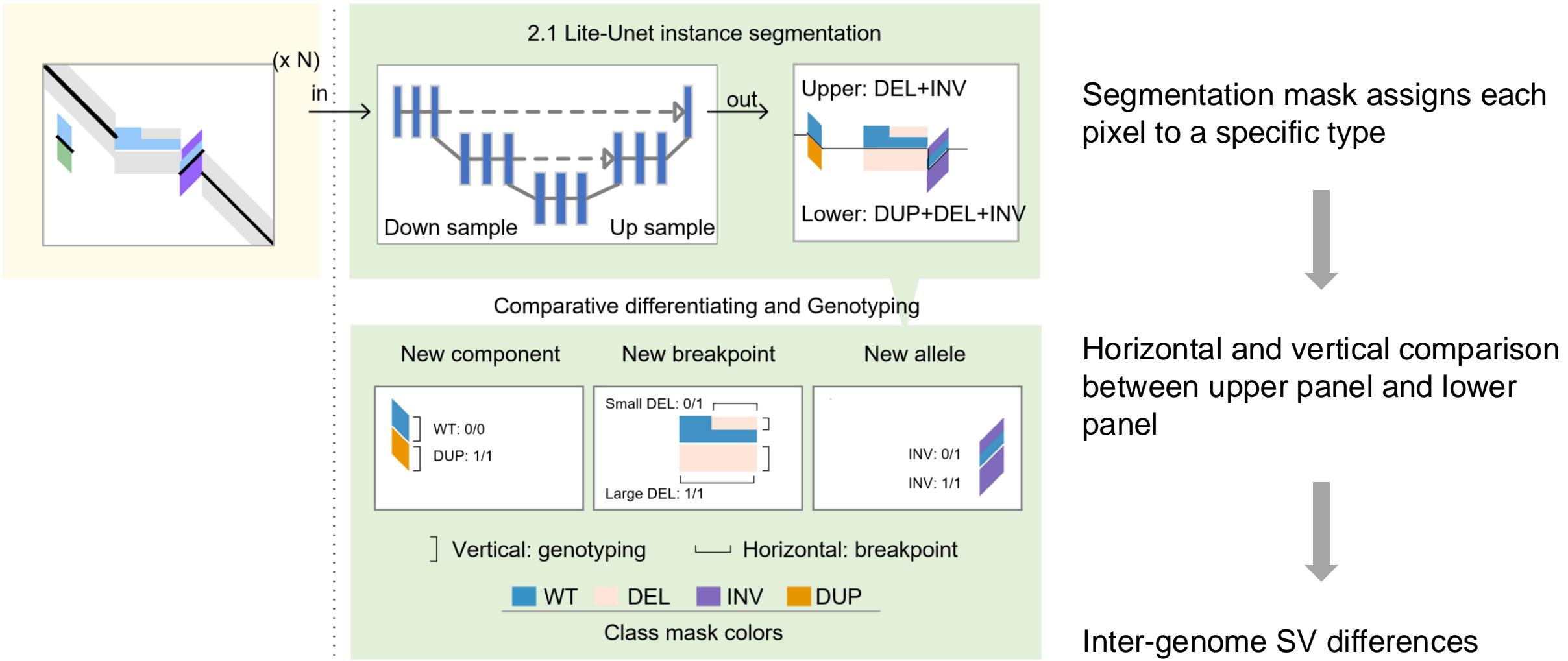
Model selection



Comparative recognition module

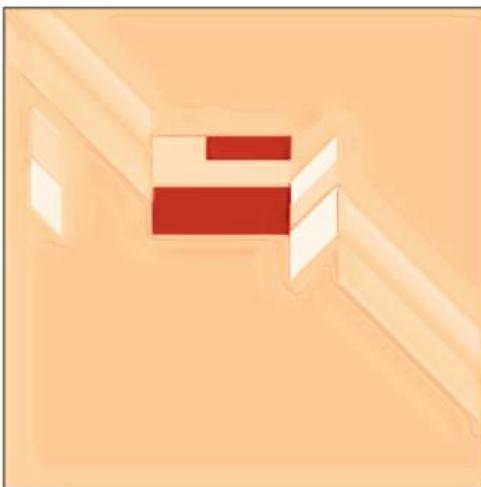


Comparative recognition module

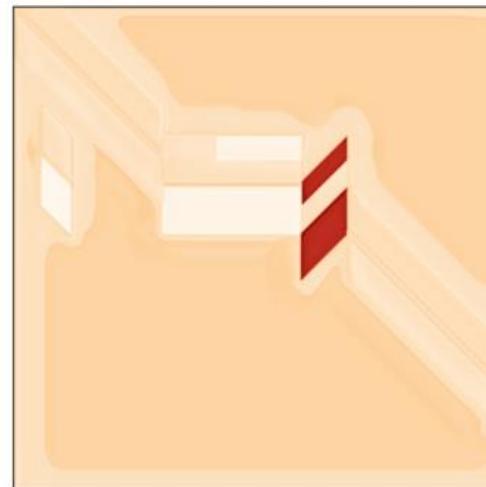


Interpretability

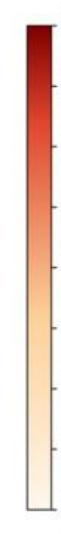
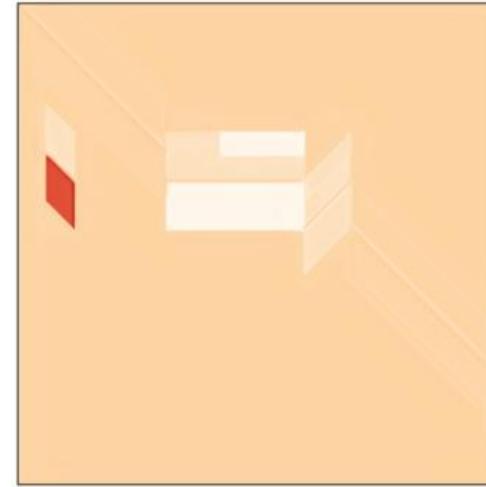
Attribution for DEL



Attribution for INV



Attribution for DUP



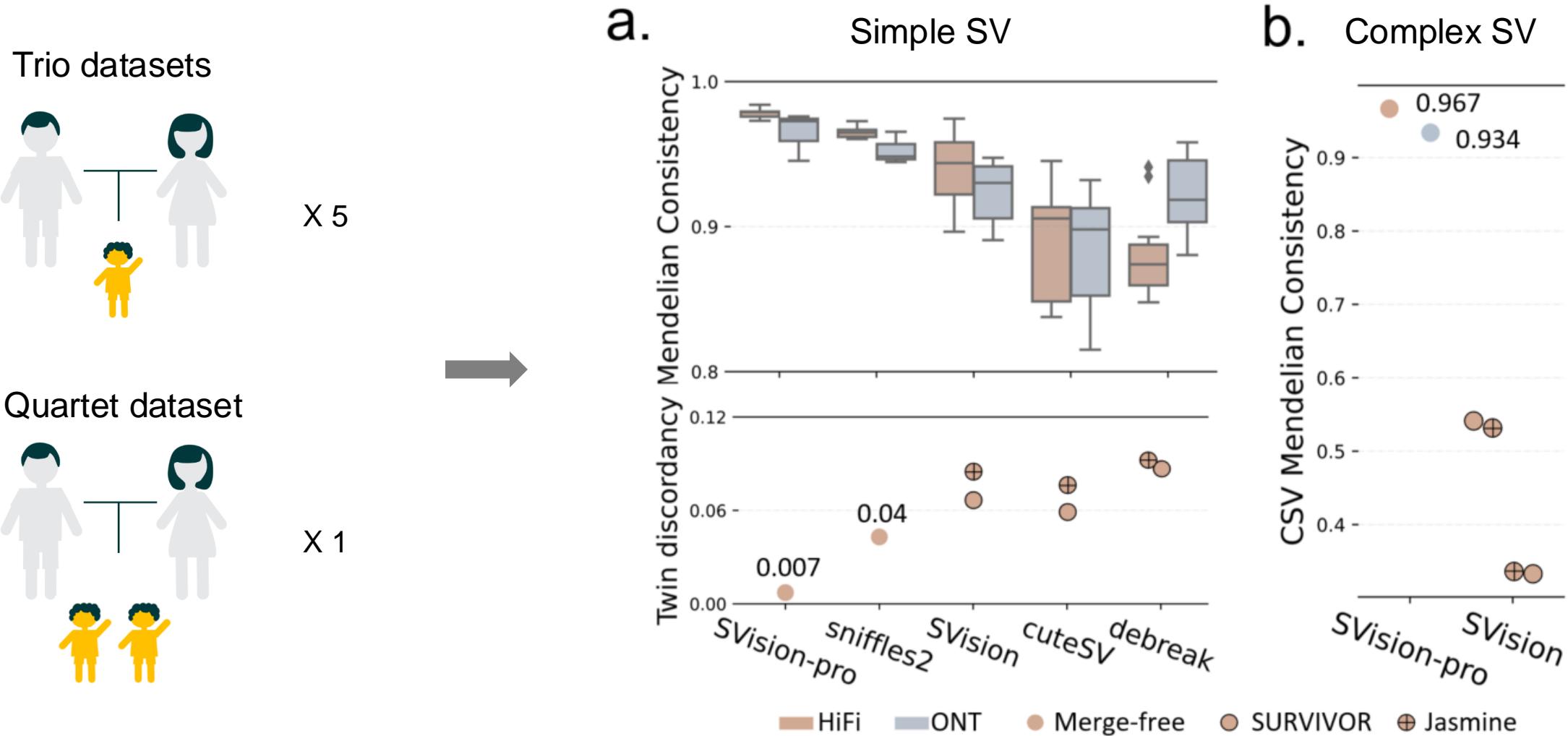
>0: Positive attribution

=0: None attribution

<0: Negative attribution

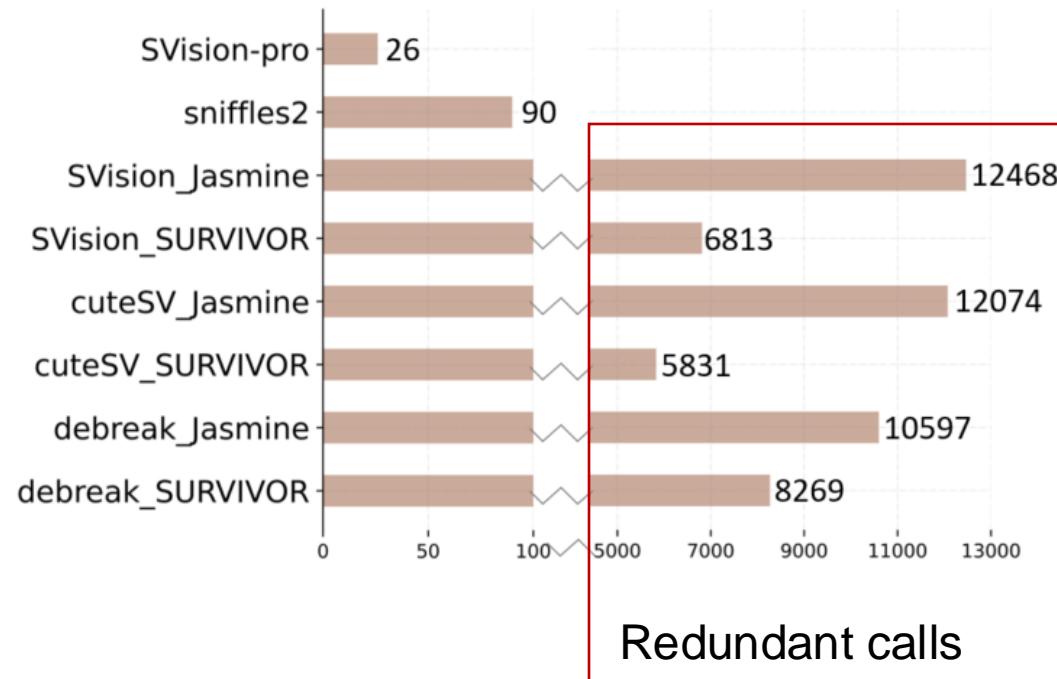
(Layer Grad-Cam)

Performance 1. Multi-genome genotyping

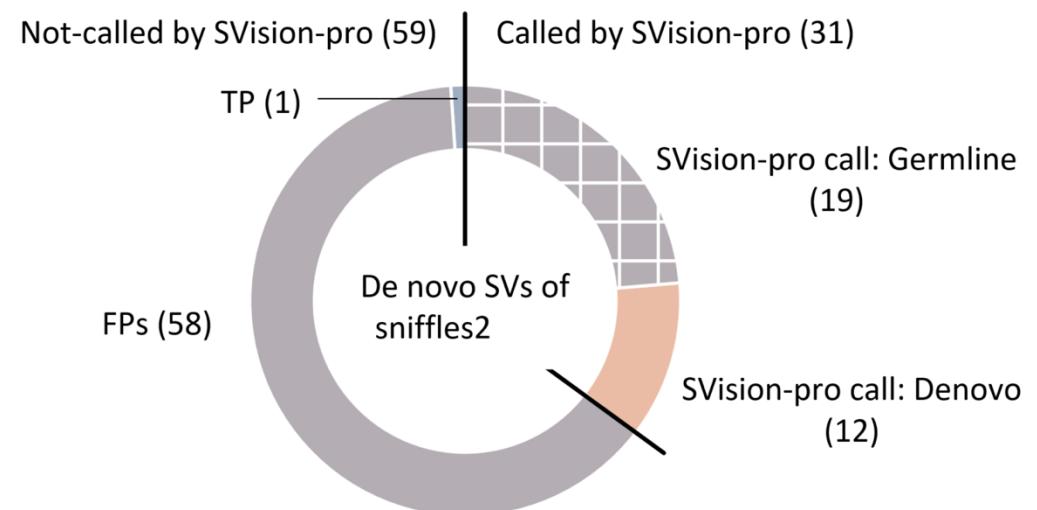


Performance 2. *De novo* SV

Number of *de novo* SVs detected from six families



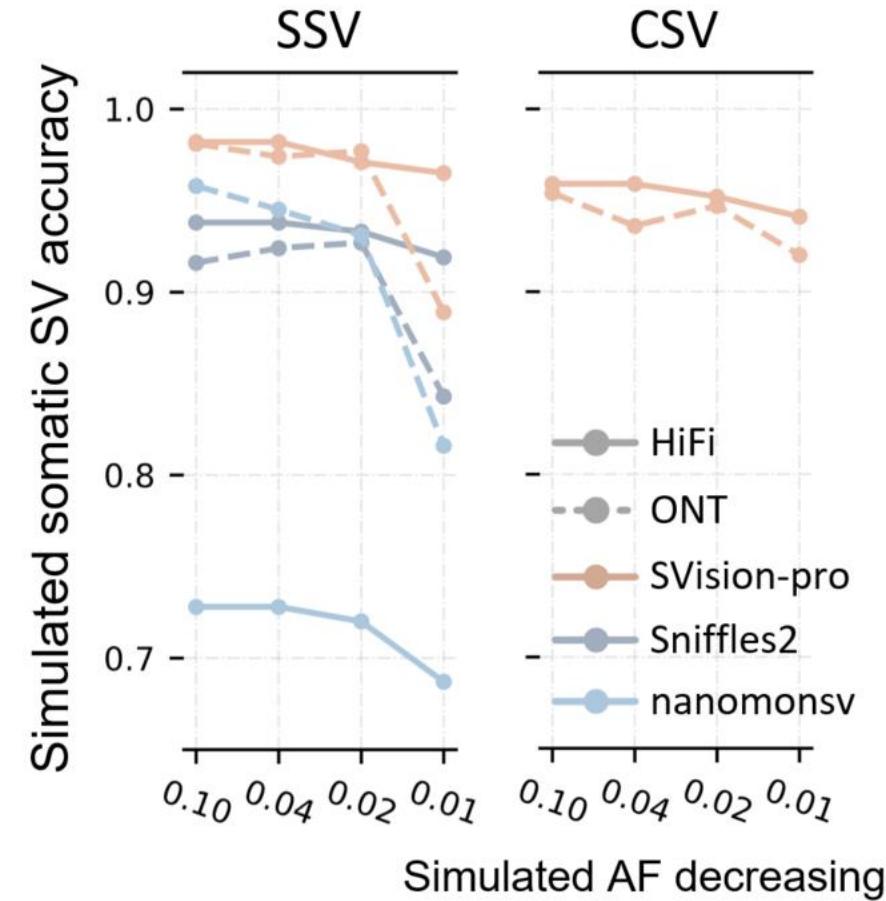
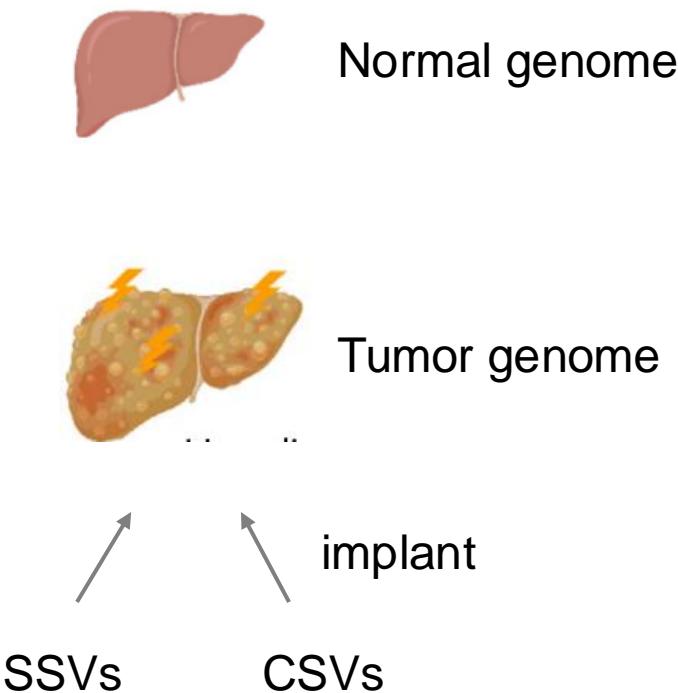
Analysis of 90 Sniffles2 *de novo* calls



SVision-pro significantly reduces false-positive *de novo* calls

Performance 3. Low-frequency somatic SV

Simulated normal-tumor pairs





Brief Communication

<https://doi.org/10.1038/s41587-024-02190-7>

De novo and somatic structural variant discovery with SVision-pro

Received: 1 August 2023

Accepted: 27 February 2024

Published online: 22 March 2024

**Songbo Wang^{1,2,3}, Jiadong Lin^{2,3}, Peng Jia^{1,2,3}, Tun Xu^{1,2,3}, Xiujuan Li^{2,3},
Yuezhuangnan Liu⁴, Dan Xu⁴, Stephen J. Bush^{2,3}, Deyu Meng^{3,5,6,7} &
Kai Ye^{1,2,3,4,8,9}✉**

Acknowledgements



西安交通大学

XIAN JIAOTONG UNIVERSITY



Mr Songbo Wang
PhD student
Automation department



Dr Jiadong Lin
Now Postdoc @ Eichler lab

Contact: kaiye@xjtu.edu.cn



中华人民共和国科学技术部

Ministry of Science and Technology of the People's Republic of China

