# VARiD: A Variation Detection Framework for Color- and Letter-Space platforms

Adrian Dalca[1,2], Steven Rumble[3], Sam Levy[4], Michael Brudno[2,5]

1. Massachusetts Institute of Technology
2. University of Toronto
3. Stanford University
4. The Scripps Research Institute
5. Donnelly Centre and the Banting and Best Department of Medical Research

# Variation detection from NGS reads

- Determine differences (variation) between **reference** and **donor** using NGS reads of the donor

```
Reference: TCAGCATCGGCATCGACTGCACAGGACCAGTCGATCGAC

Donor:         ???????????????????????????????????????
                         GCATCGACTGCA
                         CGGGATCGACTG
Aligned reads:              ATCCATTGCA
                         GATCCACTGCAC
```

# Motivation
**Color-space** and **Letter-space** platforms
bring them **together**

## Methods

## Results

## Summary

# Sequencing Platforms

- letter-space
  Sanger, 454, Illumina, etc

```
> NC_005109.2 | BRCA1 SX3
TCAGCATCGGCATCGACTGCACAGG
```

- color-space
  AB SOLiD
  less software tools available

```
> NC_005109.2 | BRCA1 AF3
T21231323031323212131112 0
```

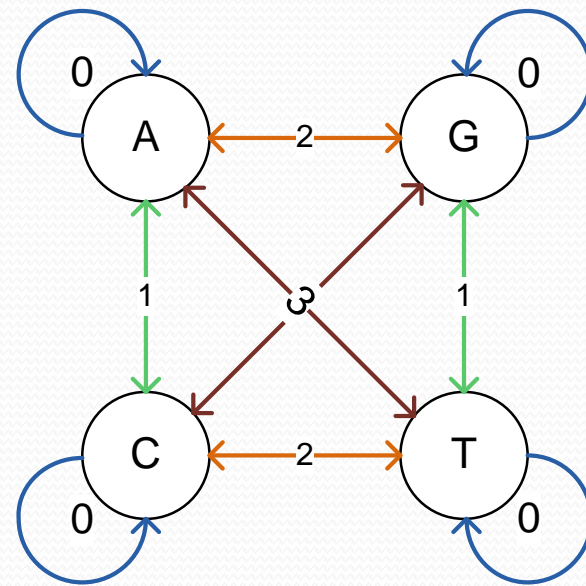- many differences -> useful to combine this information
  - sequencing biases
  - inherent errors
  - advantages

# Color Space

### Translation Matrix

|   | A | C | G | T |
|---|---|---|---|---|
| A | 0 | 1 | 2 | 3 |
| C | 1 | 0 | 3 | 2 |
| G | 2 | 3 | 0 | 1 |
| T | 3 | 2 | 1 | 0 |

### Translation Automata



> **T21231323031323212121311120**

# Color Space

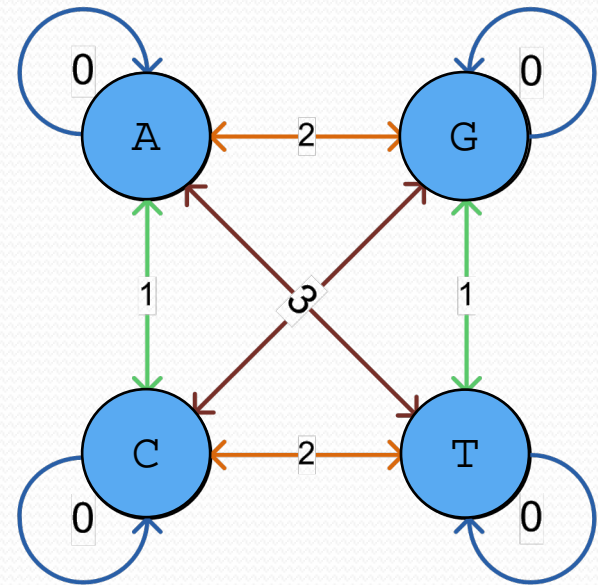Translating

> T212313230313232121311120
> TCAGCATCGGCATCGACTGCACAGG

Sequencing Error vs SNP

```
Sequencing Error
> T212313230313232121311120
> T2123132303103232121311120
> TCAGCATCGGCAAGCTGACGTGTCC

SNP
> TCAGCATCGGCATCGACTGCACAGG
> TCAGCATCGGCAGCGACTGCACAGG
> T212313230312332121311120
```
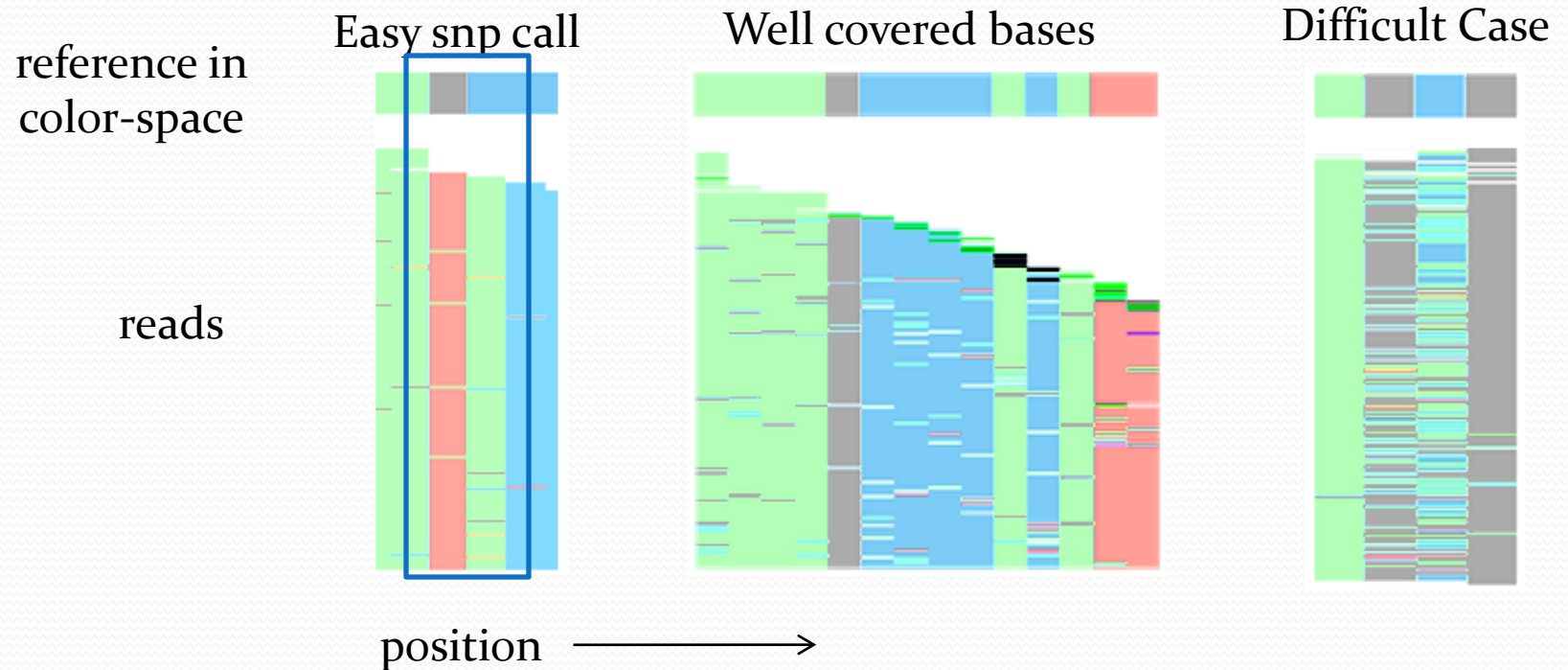
# Color Space

- clear distinction between a sequencing error and a SNP
- can this help us in SNP detection? sounds like it!
  single color change → error,
  2 colors changed → (likely) SNP.

reference in
color-space

Easy snp call

Well covered bases

Difficult Case

reads

position ⟶

# Motivation

Motivation
- variation caller  to handle both letter-space & color-space reads

Detection
- Heterozygous SNPs
- Homozygous SNPs
- Tri-allelic SNPs
- small indels
- account for various errors, quality values & misalignments

VARiD
- system to make inferences on the donor bases
  - variation detection

Motivation

# Methods

## Simple HMM Model
**states, emissions, transitions, FB**

## Extended HMM Model
gaps, diploids, exceptions
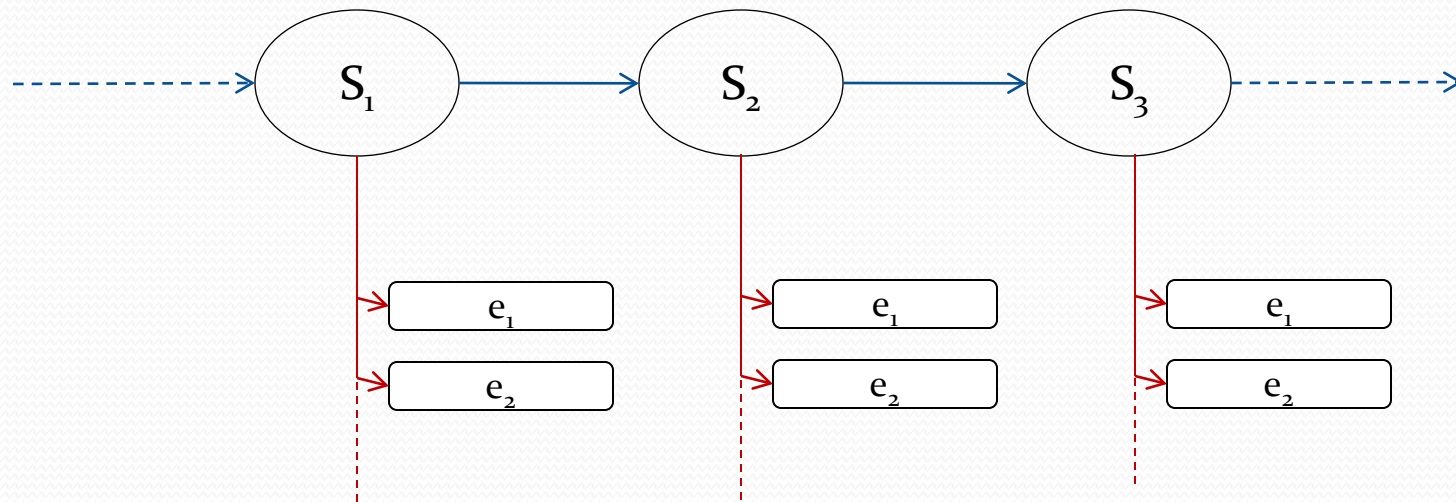
Results

Summary

# Hidden Markov Model (HMM)

**Statistical** model for a system - **states**

Assume that system is a **Markov process** with state unobserved.

Markov Process: next state depends only on current state
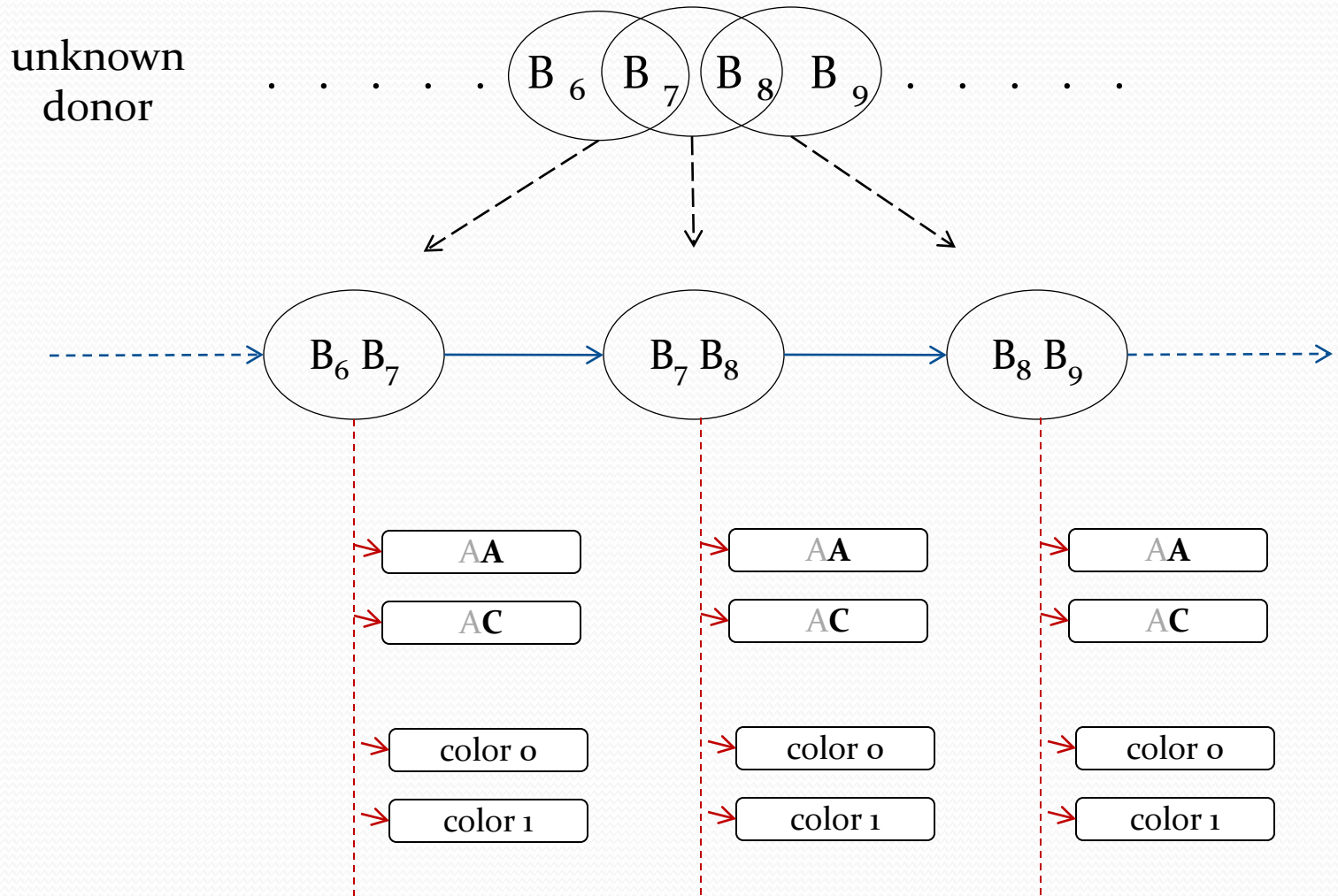
We can **observe** the state's emission (output)

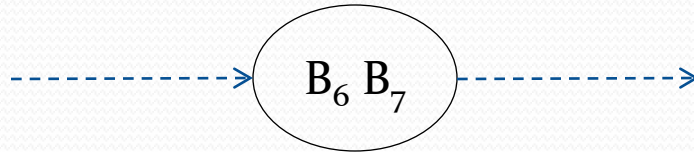each state has a probability distribution over outputs

# Hidden Markov Model (HMM)

Apply HMM to variation detection:
- we don't know the state (donor), but
- we can observe some output determined by the state (aligned reads)

# Hidden Markov Model (HMM)
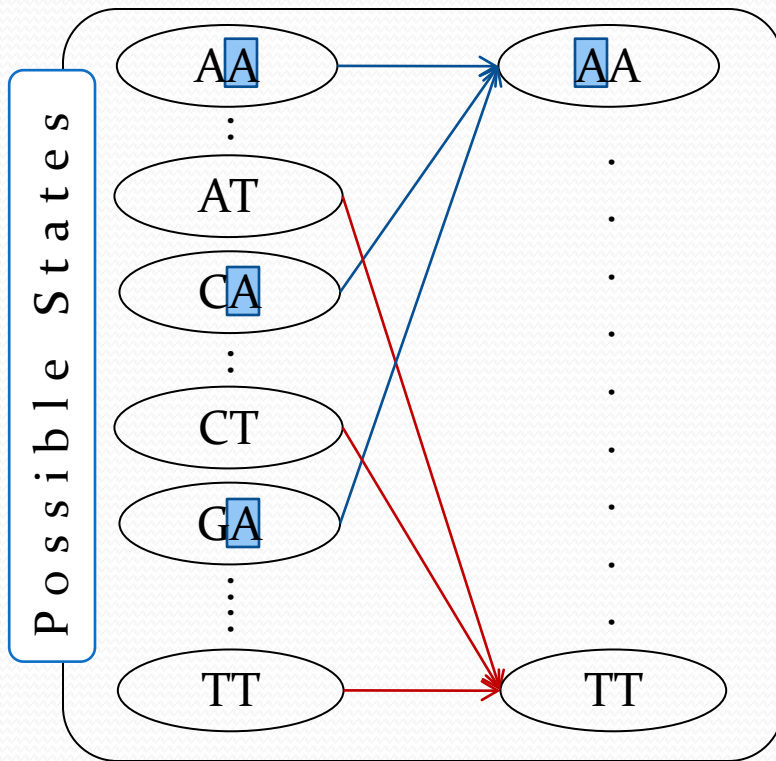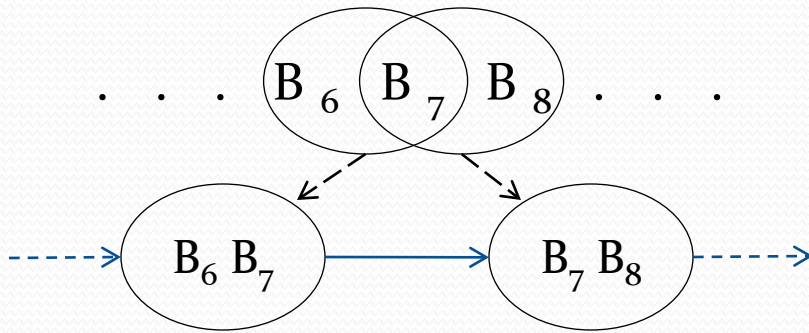
# States

$$B_6 \ B_7$$

The donor could be:
- letters: AA    color 0
- letters: AC    color 1

    :

- letters: TT    color 0

16 combinations

|   | A | C | G | T |
|---|---|---|---|---|
| A | 0 | 1 | 2 | 3 |
| C | 1 | 0 | 3 | 2 |
| G | 2 | 3 | 0 | 1 |
| T | 3 | 2 | 1 | 0 |

Why **pairs of letters**? Handle colors.
- AA and TT gives the same colors. Can't just model colors

# Transitions



## States

- 16 possible states
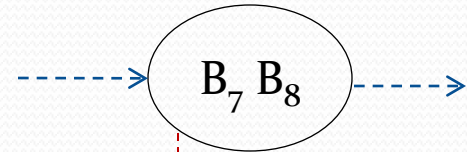- only look at second letter

## Transitions

- only certain transitions allowed
- when allowed, $p(X_t | X_{t-1}) = \text{freq}(X_t)$
- each state depends only on the previous states (Markov Process)

# Emissions

# Emission Probabilities


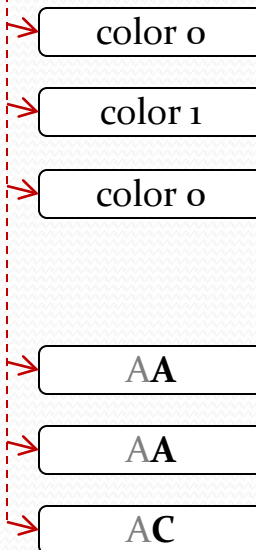
| | emission | probability p(em\|AA) | |
|---|---|---|---|
| A**A** | | | |
| | color 0 | $1 - 3\varepsilon$ | T**T** |
| | color 1 | $\varepsilon$ | Same color emission distribution |
| | color 2 | $\varepsilon$ | |
| | color 3 | $\varepsilon$ | |
| | letters A | $1 - 3\xi$ | T**T** |
| | letters C | $\xi$ | Different letter emission distribution |
| | letters G | $\xi$ | |
| | letters T | $\xi$ | |

# Emission Probabilities

$B_6 B_7 B_8 B_9$ . . . . .

Combining emission probabilities
• probability that this state emitted these reads.

```
T 0 1 0 2 0 1 0 0 3 1 1 2 2 3
 T 1 0 3 0 1 0 1 3 1 1 2 2 3
  T 2 0 1 0 0 3 1 1 2 2 3
```

E.g. For state CC:

$$p_E = [(1 - 3\varepsilon)^2 \times \varepsilon^1] \times [(1 - 3\xi)^1 \times \xi^2]$$

```
A T T G C G C A A T G C G
 T T G G G C A A T G C G A
  G C G C A C T G C G A C
```
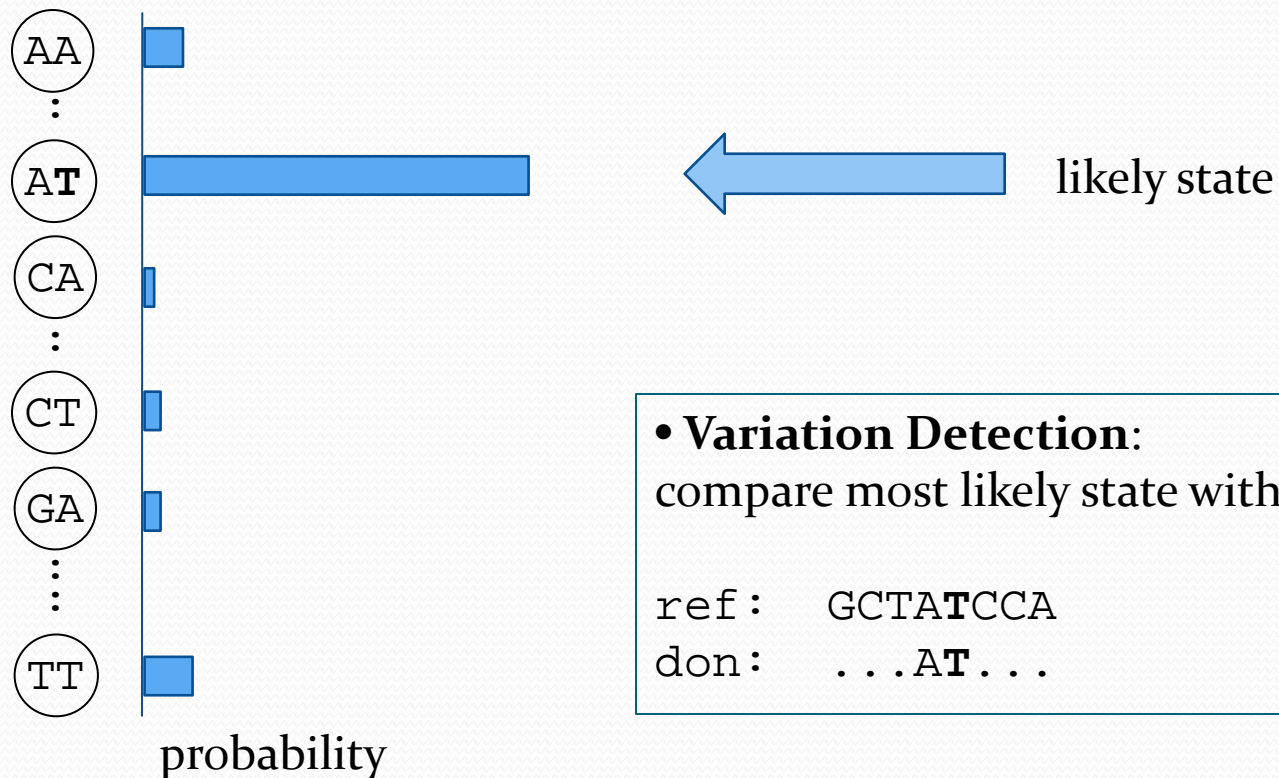
# Simple HMM

Summary

- unknown state
  - donor pair at location

- transitions
  - transition probabilities

- emissions
  - reads at location
  - emission probabilities

$B_6$ $B_7$

AA

AC

color 0

color 1

# Forward-Backward Algorithm

- Have set-up a form of an HMM
- run Forward-Backward algorithm
- get probability distribution over states at each position



likely state

- **Variation Detection**:
compare most likely state with reference:

```
ref:   GCTATCCA
don:   ...AT...
```

probability

Motivation

# Methods

## Simple HMM Model
states, emissions, transitions, FB

## Extended HMM Model
gaps, diploids, exceptions

Results

Summary

# Extended HMM

Simple HMM
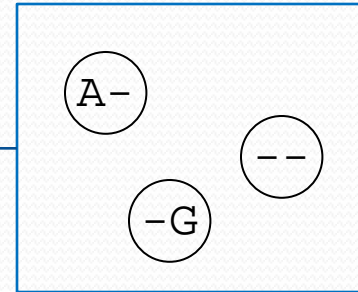- only detects **homozygous SNPs**

Extended HMM:
- short indels
- heterozygous SNPs
- complex error profiles & quality values

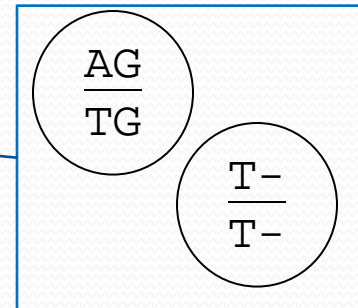# Expansion: Gaps and heterozygous SNPs

Expand states
- Have states that include gaps
  - emit: gap or color



- Have larger states, for diploids



- Transitions built in similar fashion as before
- Same algorithm, but in all we have 1600 states with very sparse transitions

# Expansion

- Emission probabilities
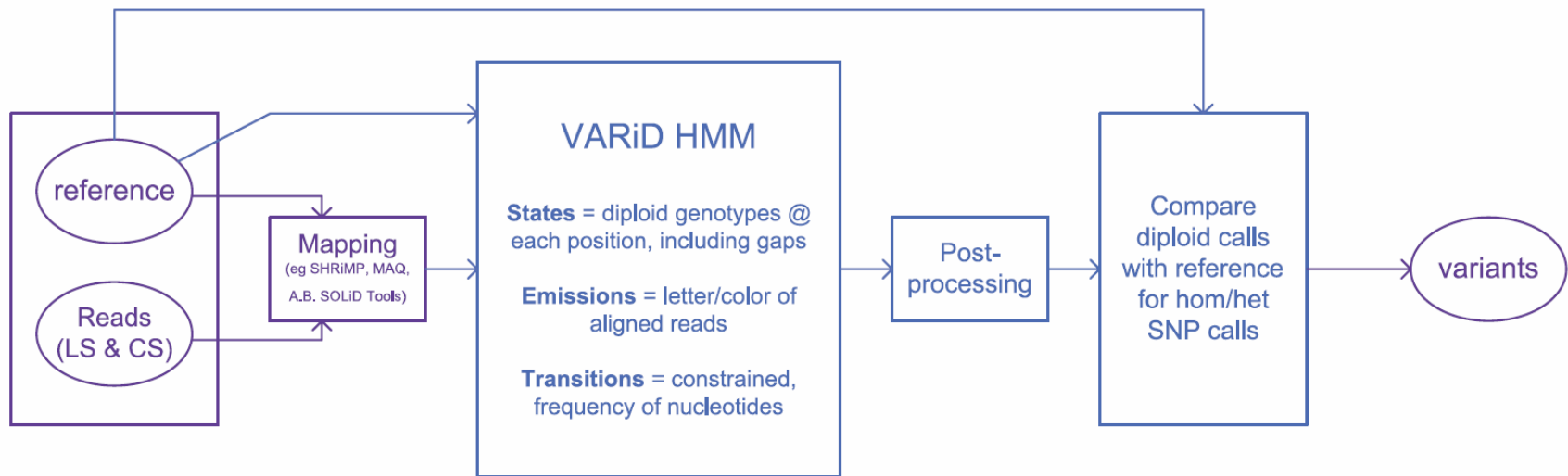  - Support quality values
  - Use variable error rates for emissions

- Translate through the first color
  - first color is incorrect
  - letter-space signal

```
Donor:  ACAGCATCGGCATCGACTGC
        112313230312332 1213
read:  >T2112313230312332 1213
       > C112313230312332 1213
```

- Post-process putative SNPs
  - correlated adjacent errors may support het SNPs
  - check putative SNPs

# Summary



blue: varid steps

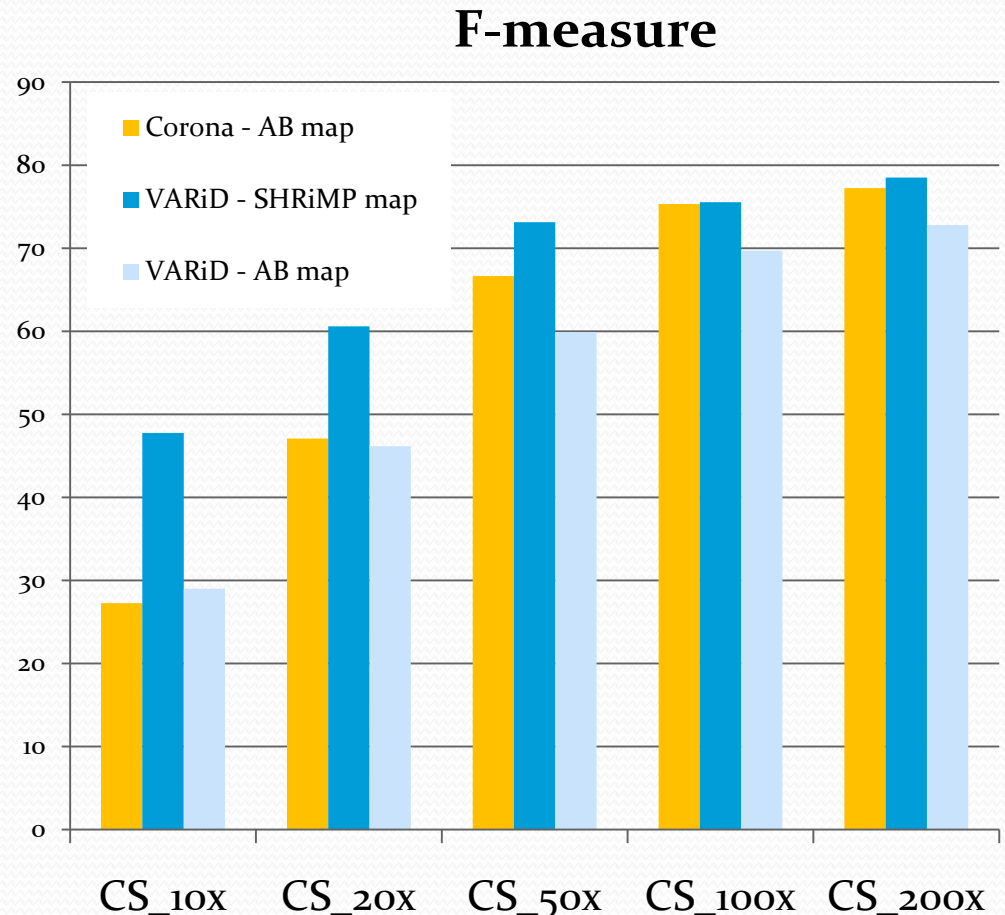Motivation

Methods

Results

Summary

# Results

• Human dataset from Harismendy et al, 2009. (NA17156,17275,17460,17773) 454, SOLiD, Sanger

**Color-space** dataset:
• Compare random subsets:
  • Corona (with AB mapper)
  • VARiD (with SHRiMP)
  • VARiD (with AB mapper)

**Conclusions**:
• the three pipelines perform very similarly.
• High-coverage results is as good as can be achieved

**F-measure**



Legend:
- Corona - AB map
- VARiD - SHRiMP map
- VARiD - AB map

x-axis: CS_10x, CS_20x, CS_50x, CS_100x, CS_200x

# Results

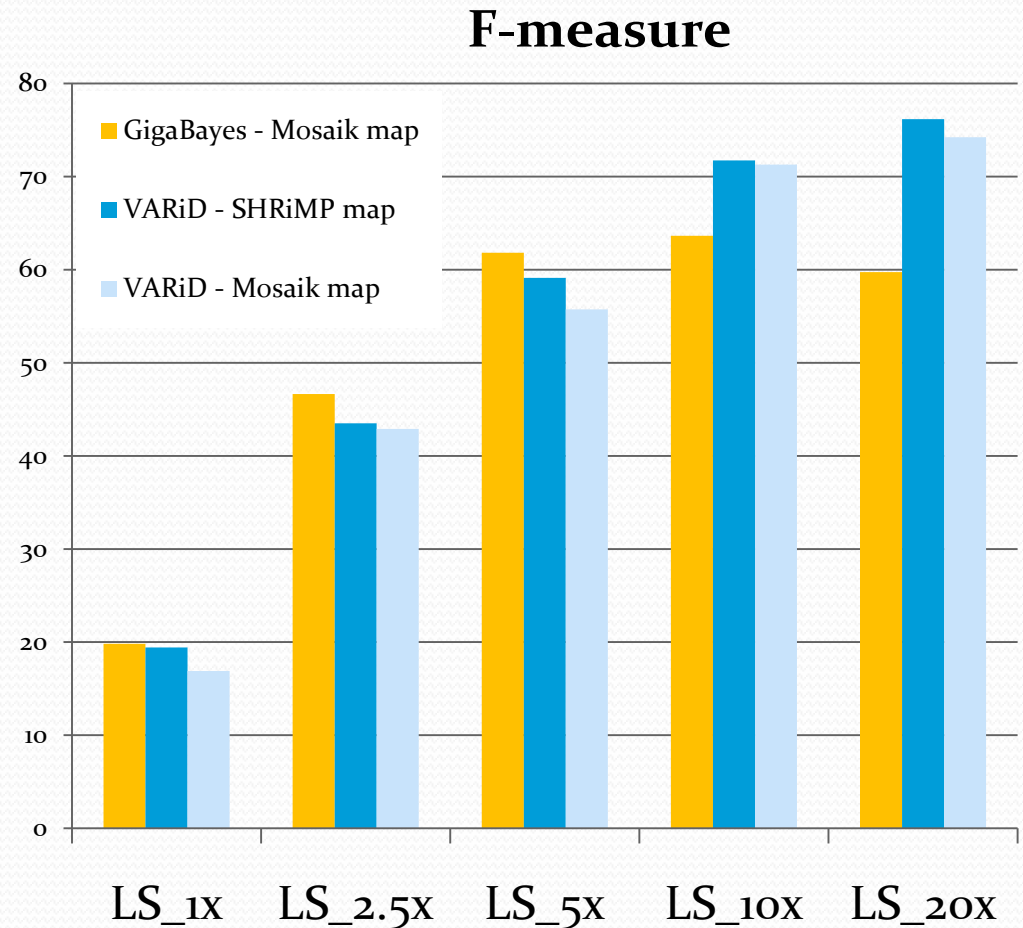**Letter-space** dataset:
- Compare random subsets :
  - GigaBayes (with Mosaik)
  - VARiD (with SHRiMP)
  - VARiD (with Mosaik)

**Conclusion**:
- the three pipelines perform very similarly.
- High-coverage results is as good as can be achieved



**F-measure**

Legend:
- GigaBayes - Mosaik map
- VARiD - SHRiMP map
- VARiD - Mosaik map

Categories: LS_1x, LS_2.5x, LS_5x, LS_10x, LS_20x

# Results

**VARiD Motivation:**
Combining Letter-space and Color-space data to achieve increased accuracy in at-cost comparison

Assuming a same-cost comparison of:

- 10x letter-space (LS)
- 100x color-space (CS)
- 5x LS and 50x CS



VARiD

|  | letter-space | | | | |
|---|---|---|---|---|---|
| **F-meas.** | 0x | 1x | 2.5x | 5x | 10x |
| 0x | 0.0 | 19.4 | 43.5 | 59.1 | 71.7 |
| 10x | 47.8 | 51.8 | 59.4 | 69.5 | 76.5 |
| 20x | 60.6 | 58.9 | 65.3 | 73.4 | 80.3 |
| 50x | 73.1 | 69.8 | 73.6 | 80.0 | 83.5 |
| 100x | 75.6 | 75.2 | 77.9 | 82.7 | 86.0 |

Color-space

Motivation

Methods

Results

Summary

# Summary

**Summary of VARiD**

- HMM modeling underlying donor

- Treats color-space and letter-space together in the same framework

- no translation – take advantage of each technology's properties

- accurately calls SNPs, **short indels** in both color- and letter-space

  - improved results with hybrid data.

- **Website**: http://compbio.cs.utoronto.ca/varid
(VARiD freely available)

- **Contact**: varid@cs.utoronto.ca

# Acknowledgements

**Acknowledgements**
- Steven Rumble, Sam Levy, Michael Brudno
- Misko Dzamba

**Funding**



**Partial Participant Support Award** :

- ISCB, Department of Energy, and National Science Foundation

- **Website**: http://compbio.cs.utoronto.ca/varid
(VARiD freely available)

- **Contact**: varid@cs.utoronto.ca