

# Stochastic Models for DNA Tandem Duplication

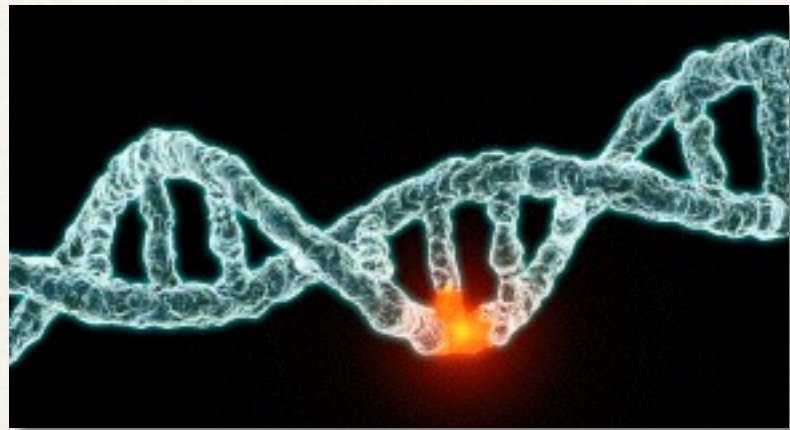
Farzad Farnoud, with M. Schwartz, J. Bruck

*Jan 18, 2016, University of Washington*



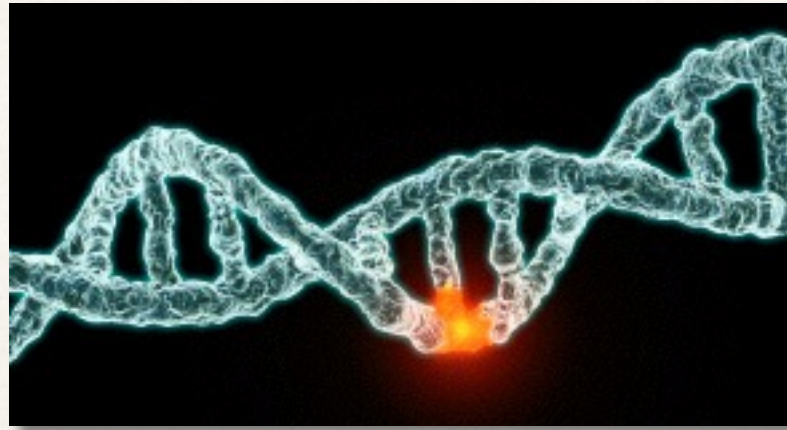


# Mutations





# Mutations

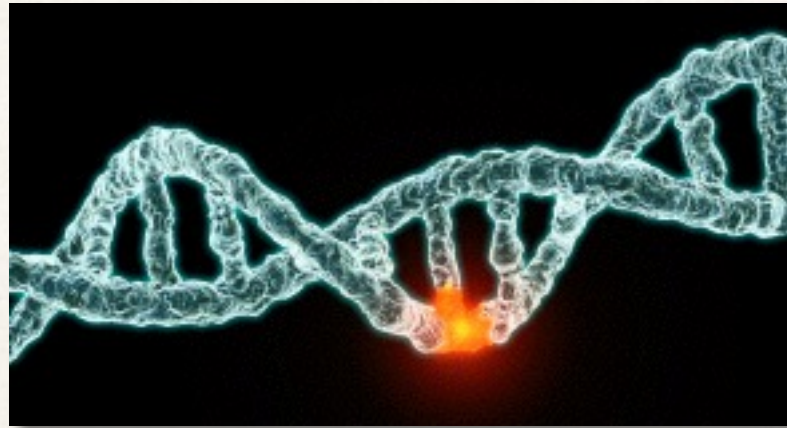


8.7 million species





# Mutations



Data storage in *live DNA*

8.7 million species





# Types of Mutations

---

TGATGCA

↓ Point Mutation

TCATGCA

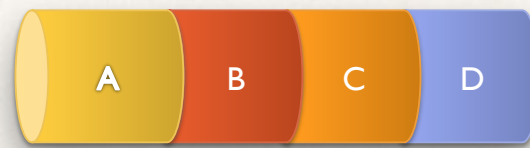
# Types of Mutations

---

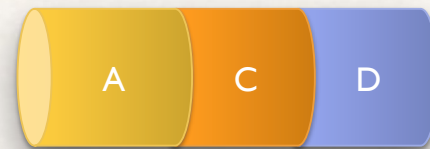
TGATGCA

↓ Point Mutation

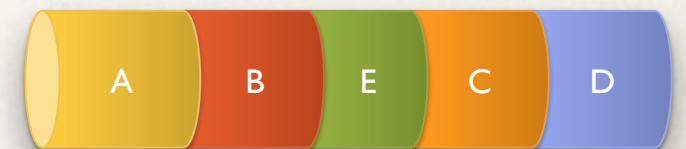
TCATGCA



↓ Deletion



↓ Insertion





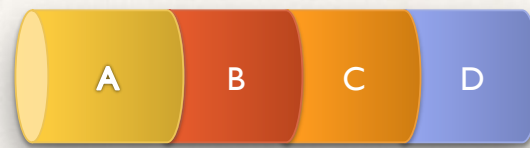
# Types of Mutations

---

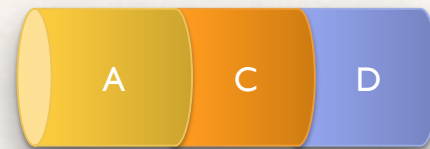
TGATGCA

↓ Point Mutation

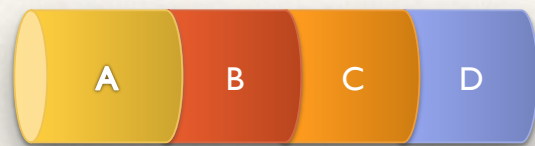
TCATGCA



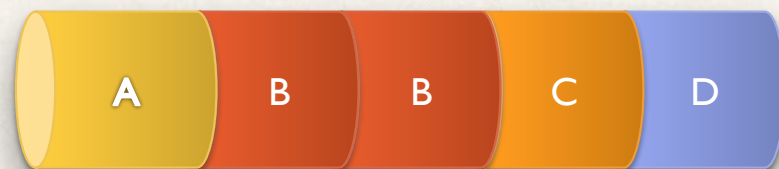
↓ Deletion



↓ Insertion



↓ Tandem Duplication



↓ Interspersed Duplication



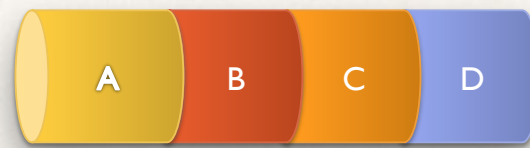
# Types of Mutations

---

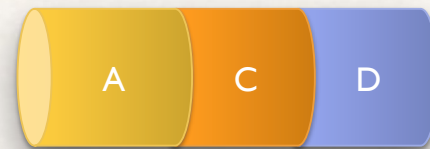
TGATGCA

↓ Point Mutation

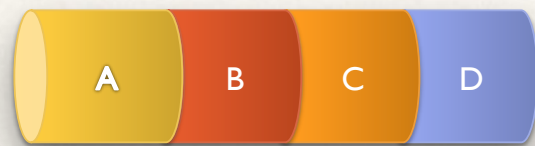
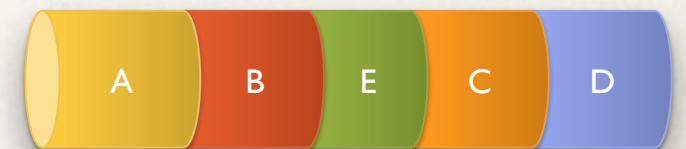
TCATGCA



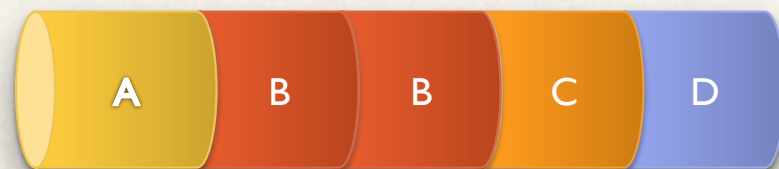
↓ Deletion



↓ Insertion



↓ Tandem Duplication



↓ Interspersed Duplication



3% of human genome



















# Point mutations are in the same positions

---

GCT <b>C</b> CGTTACAGGTGGGC <b>A</b> GGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGG <b>A</b> GGC <b>G</b>	GCT <b>C</b> CGTTACAGGTGGGC <b>A</b> GGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGG <b>A</b> GGC <b>G</b>	GCT <b>C</b> CGTTACAGGTGGGC <b>A</b> GGGG <b>A</b> GGC <b>G</b>
GCT <b>C</b> CGTTACAGGTGGGC <b>A</b> GGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGG <b>G</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>A</b> GGC <b>G</b>	GCT <b>C</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>G</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>A</b> GGC <b>T</b>	GCT <b>C</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>A</b> GGC <b>T</b>
GCT <b>C</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>G</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>G</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>A</b> GGC <b>G</b>
GCT <b>C</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>A</b> GGC <b>G</b>	GCT <b>C</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGG <b>A</b> GGC <b>G</b>



# Stochastic Model

---

Start from one repeat unit (*seed*).

Random mutations:

- Tandem duplications (TD)  
  of one or more repeat units
- Point mutations (PM)



# Stochastic Model

---

ACGT

Start from one repeat unit (*seed*).

Random mutations:

- Tandem duplications (TD)  
  of one or more repeat units
- Point mutations (PM)



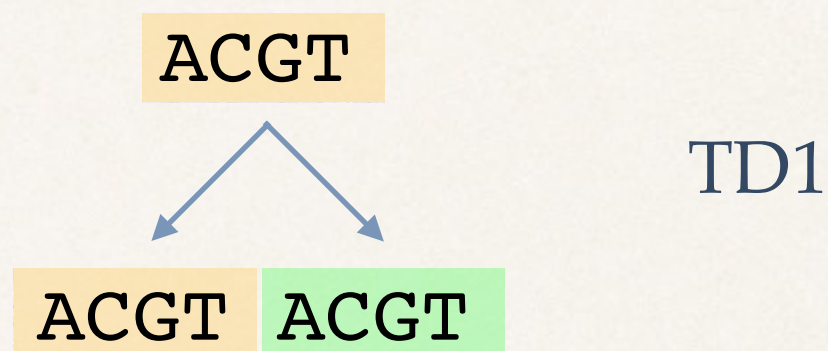
# Stochastic Model

---

Start from one repeat unit (*seed*).

Random mutations:

- Tandem duplications (TD)  
of one or more repeat units
- Point mutations (PM)





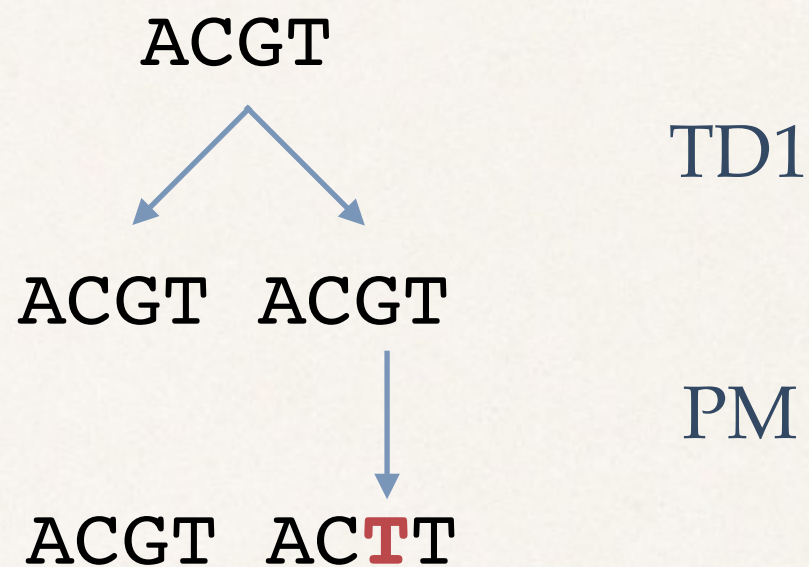
# Stochastic Model

---

Start from one repeat unit (*seed*).

Random mutations:

- Tandem duplications (TD)  
of one or more repeat units
- Point mutations (PM)





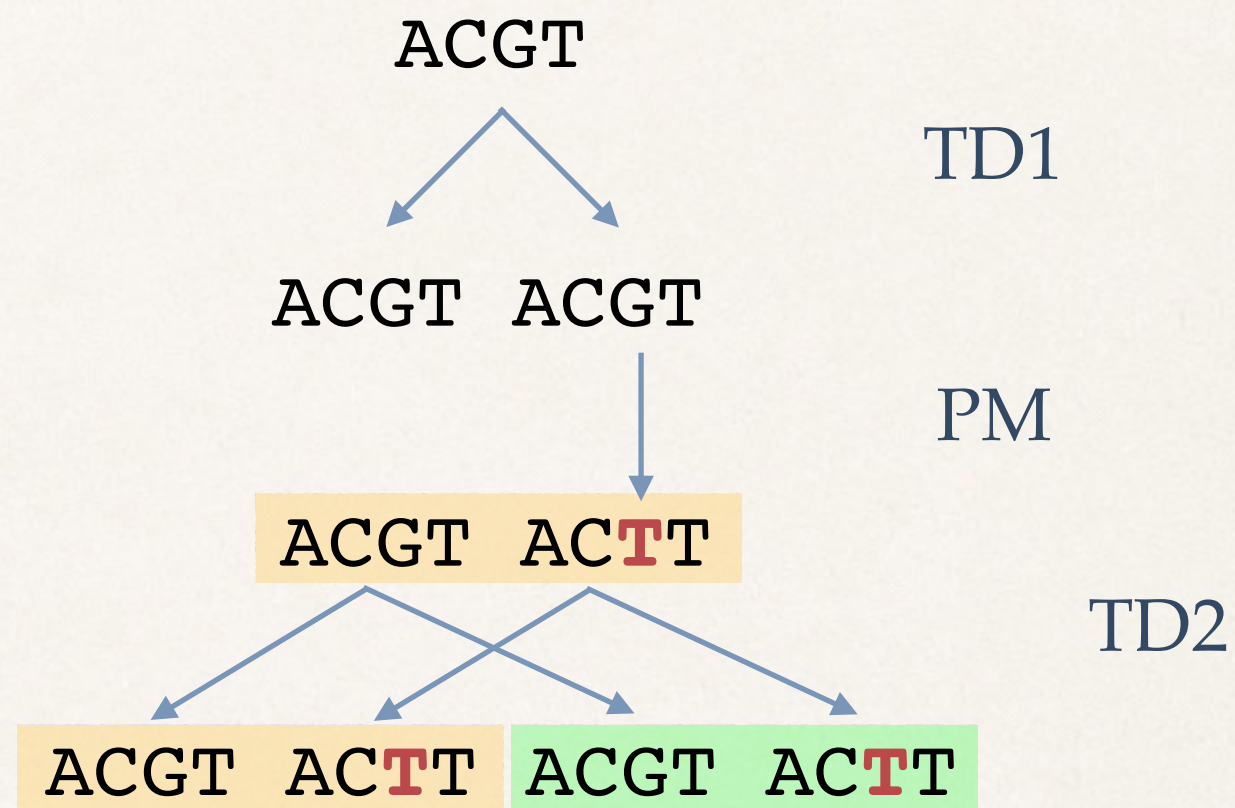
# Stochastic Model

---

Start from one repeat unit (*seed*).

Random mutations:

- Tandem duplications (TD)  
of one or more repeat units
- Point mutations (PM)





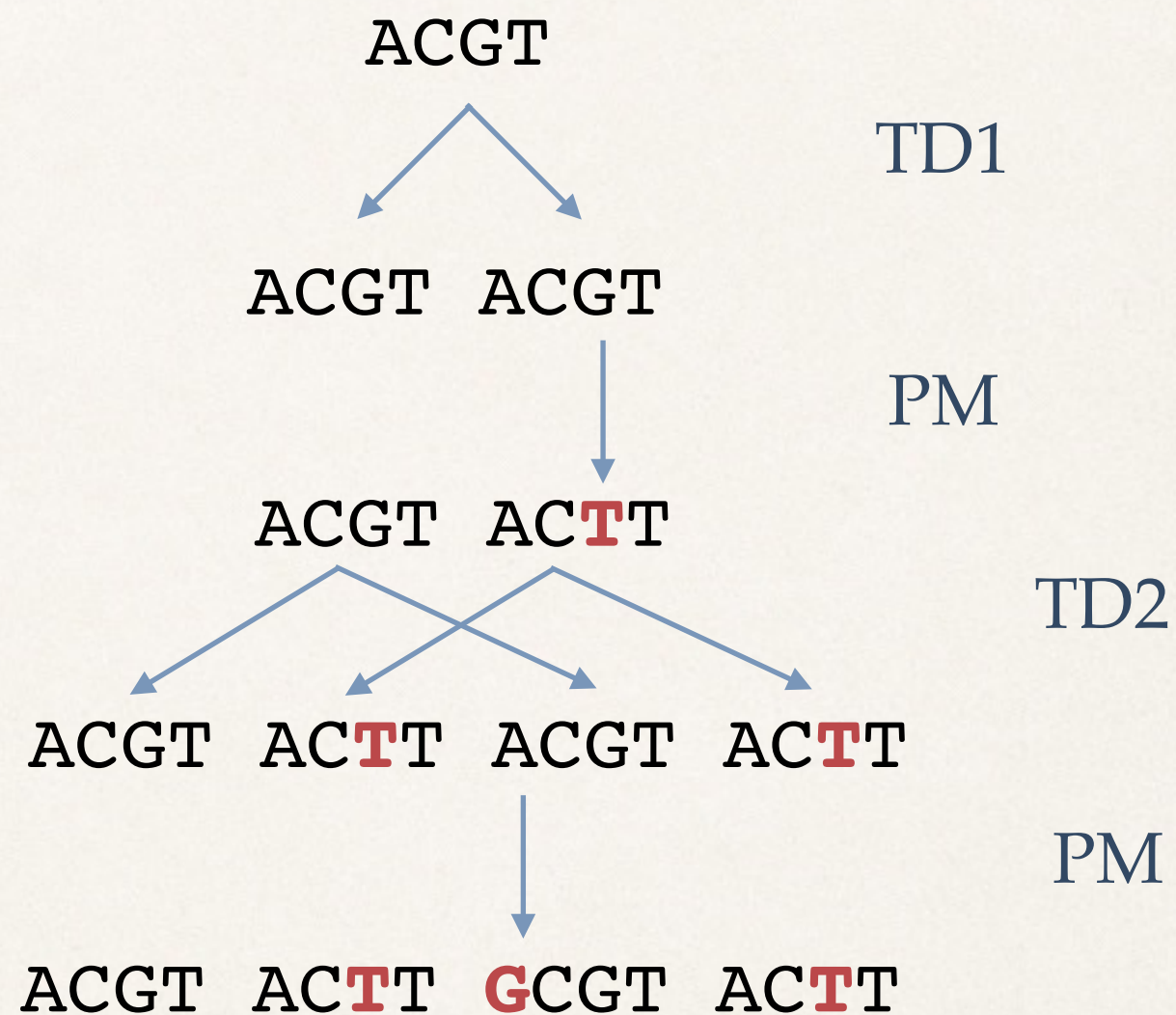
# Stochastic Model

---

Start from one repeat unit (*seed*).

Random mutations:

- Tandem duplications (TD)  
of one or more repeat units
- Point mutations (PM)



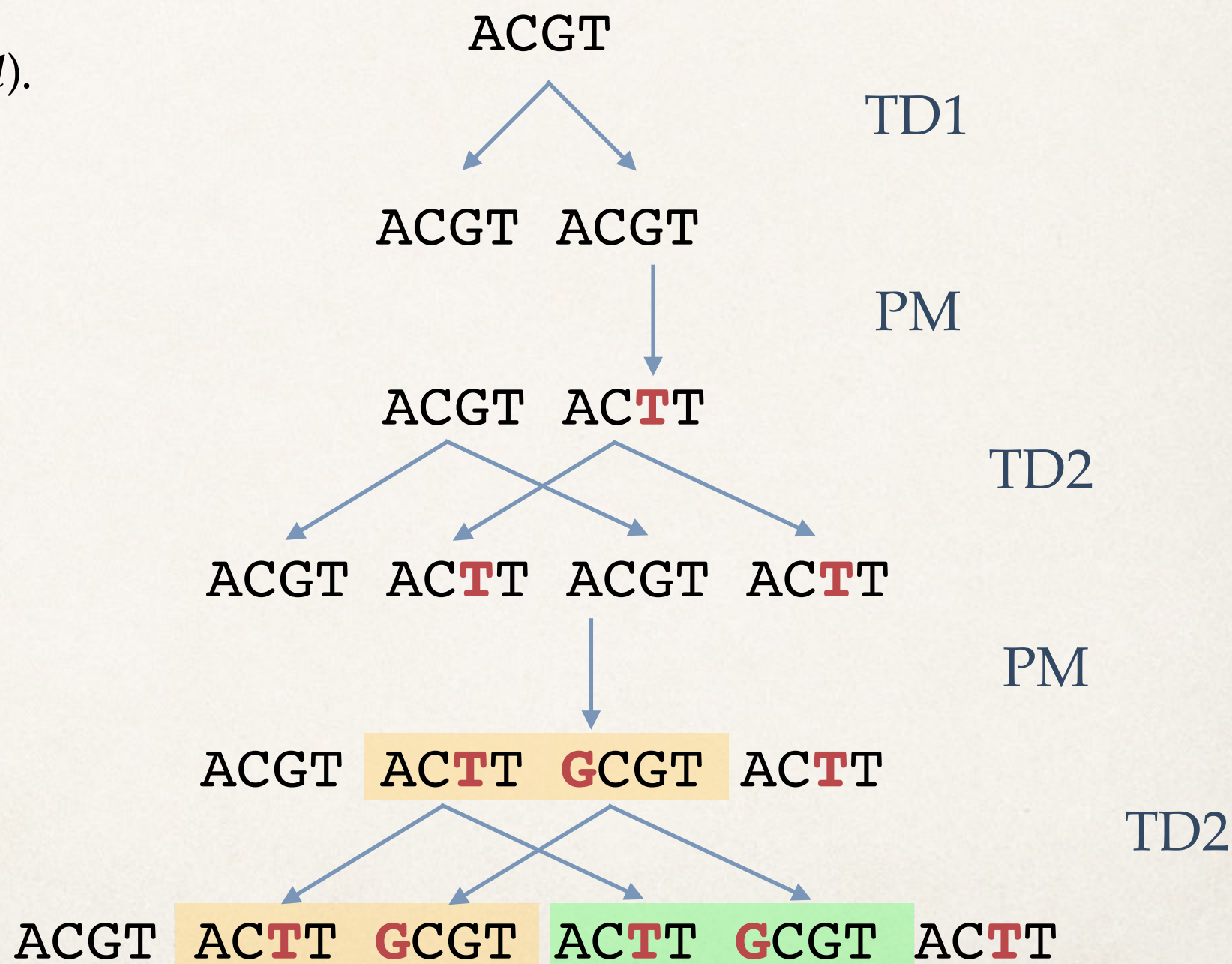


# Stochastic Model

Start from one repeat unit (*seed*).

Random mutations:

- Tandem duplications (TD)  
of one or more repeat units
- Point mutations (PM)





# Stochastic Model

Start from one repeat unit (*seed*).

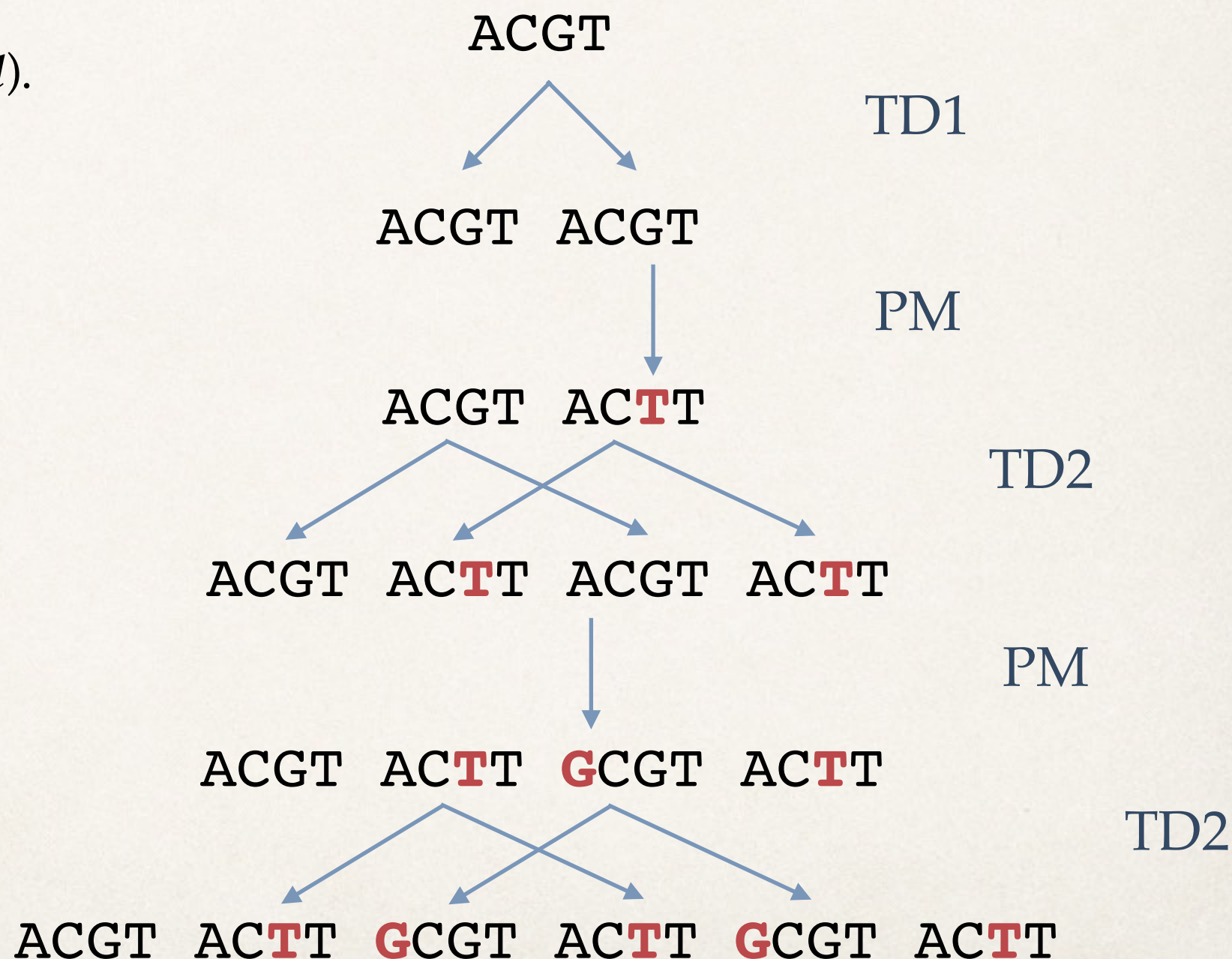
Random mutations:

- Tandem duplications (TD)  
of one or more repeat units
- Point mutations (PM)

Parameters of the model:

- Prob. of PM
- Prob. of TDs  
of different lengths

Can we learn them?





# Stochastic Model

---

Start from one repeat unit (*seed*).

Random mutations:

- Tandem duplications (TD)  
of one or more repeat units
- Point mutations (PM)

Parameters of the model:

- Prob. of PM
- Prob. of TDs  
of different lengths

Can we learn them?

ACGT ACTT GCGT ACTT GCGT ACTT



# Finding Duplication History

---

ACGT ACTT GCGT ACTT GCGT ACTT



# Finding Duplication History

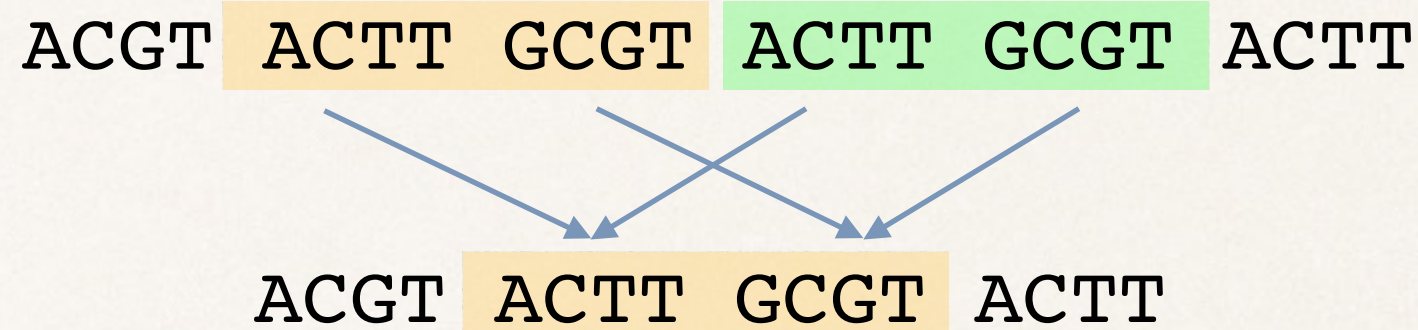
---

ACGT ACTT GCGT ACTT GCGT ACTT



# Finding Duplication History

---



TD2



# Finding Duplication History

---

ACGT ACTT GCGT ACTT GCGT ACTT

TD2

ACGT ACTT GCGT ACTT

PM

ACGT ACTT ACGT ACTT



# Finding Duplication History

---

ACGT ACTT GCGT ACTT GCGT ACTT

TD2

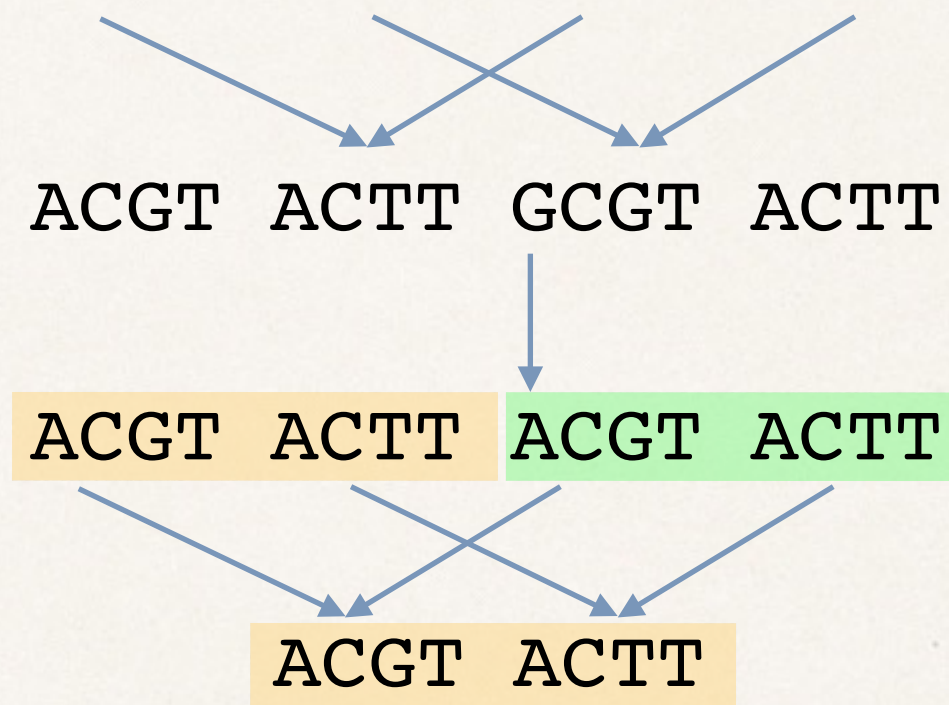
ACGT ACTT GCGT ACTT

PM

ACGT ACTT ACGT ACTT

TD2

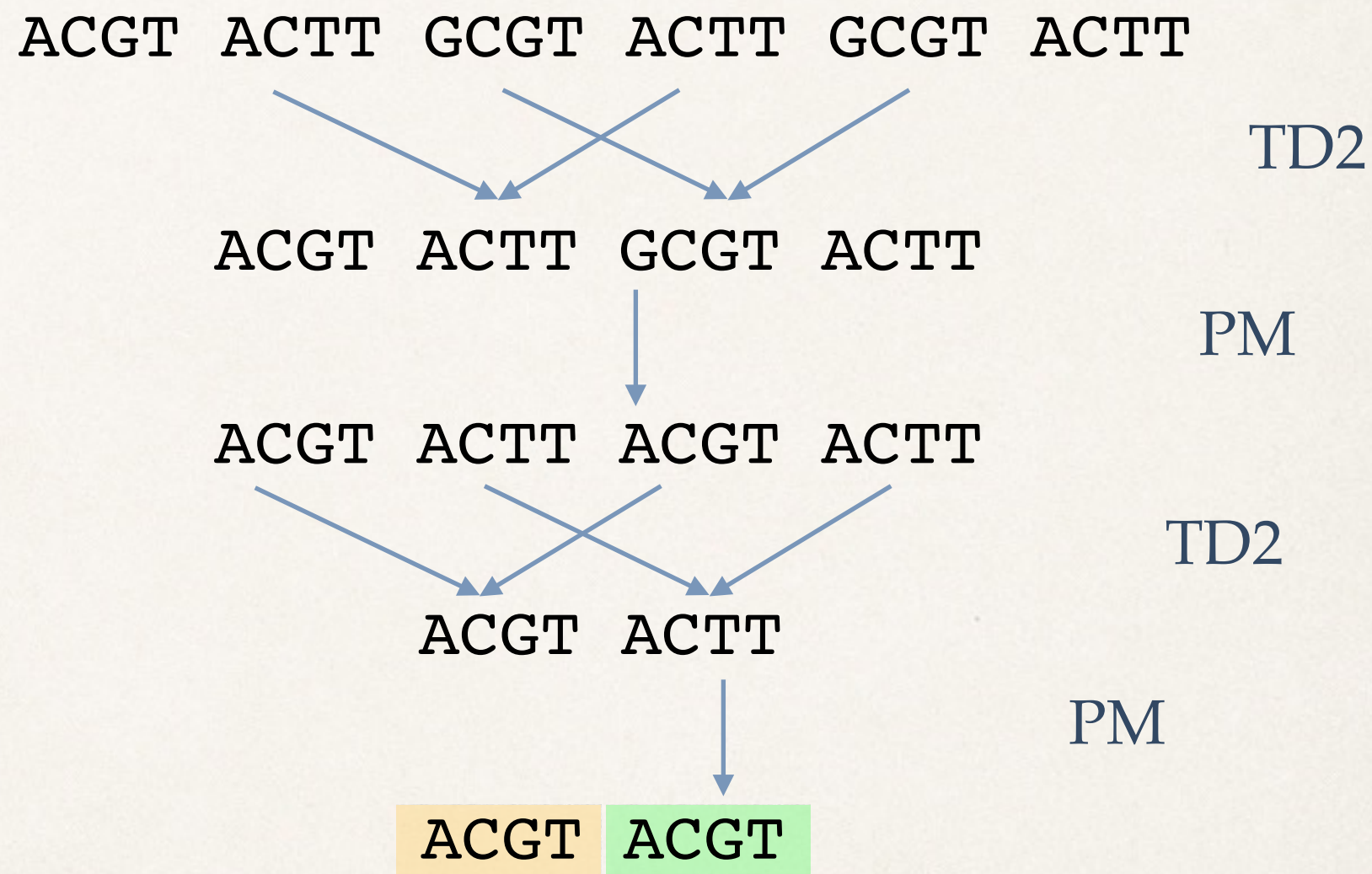
ACGT ACTT





# Finding Duplication History

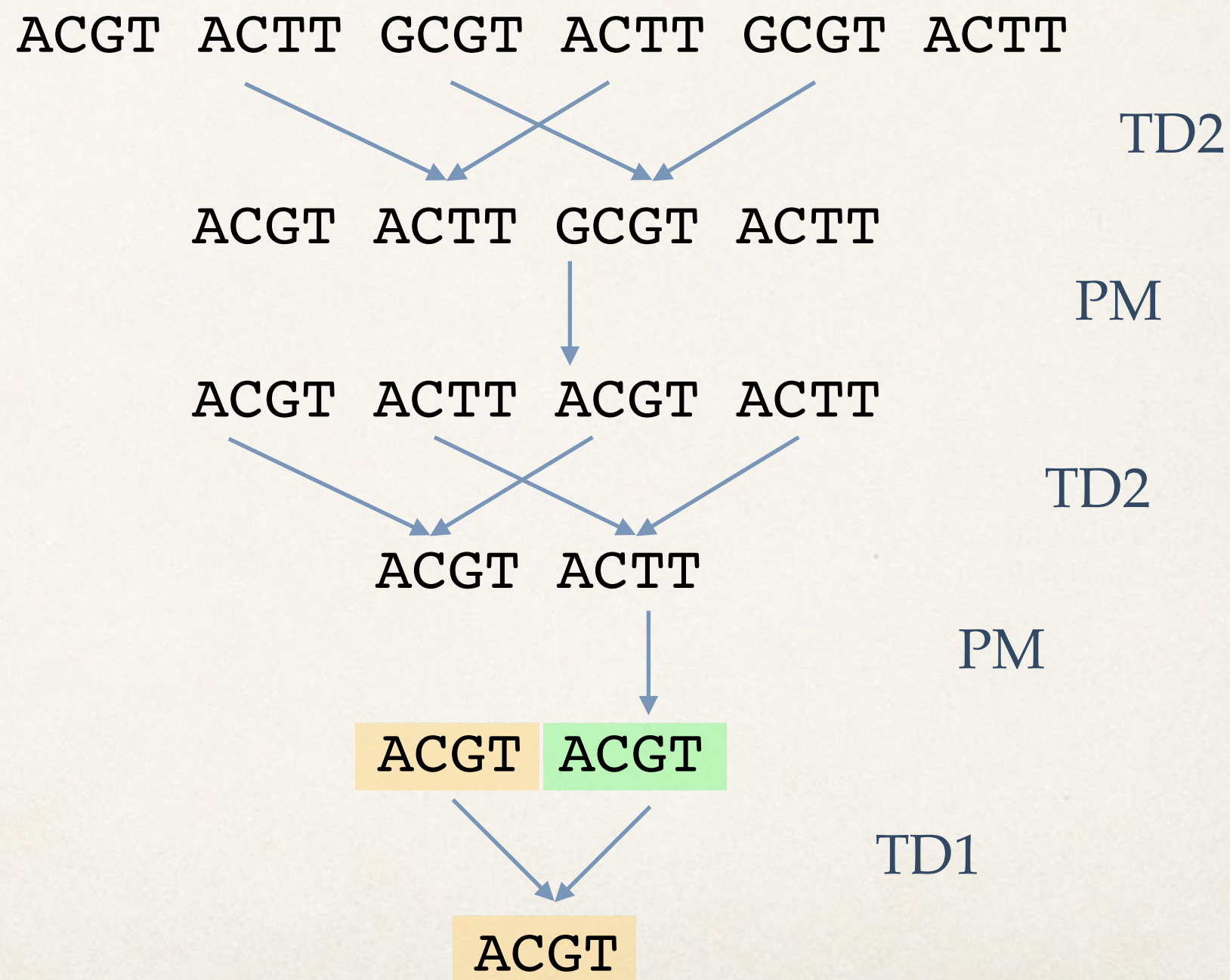
---





# Finding Duplication History

---





# Finding Duplication History

---

ACGT ACTT GCGT ACTT GCGT ACTT

TD2

ACGT ACTT GCGT ACTT

PM

ACGT ACTT ACGT ACTT

TD2

ACGT ACTT

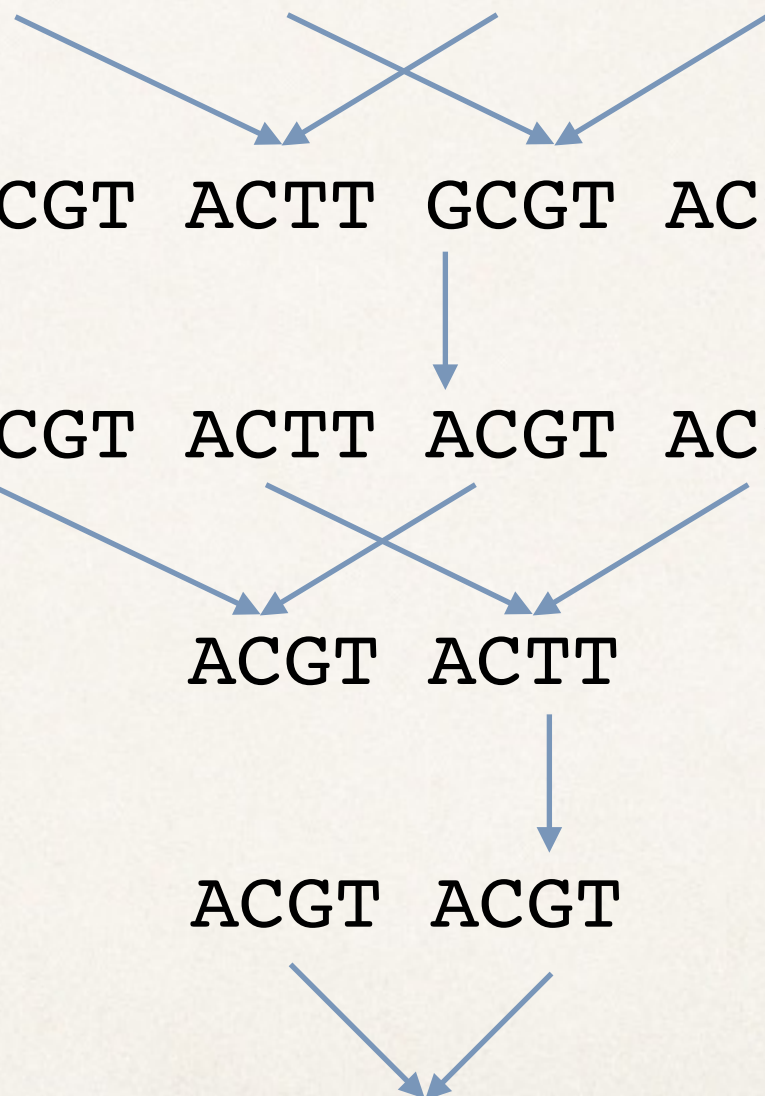
PM

ACGT ACGT

TD1

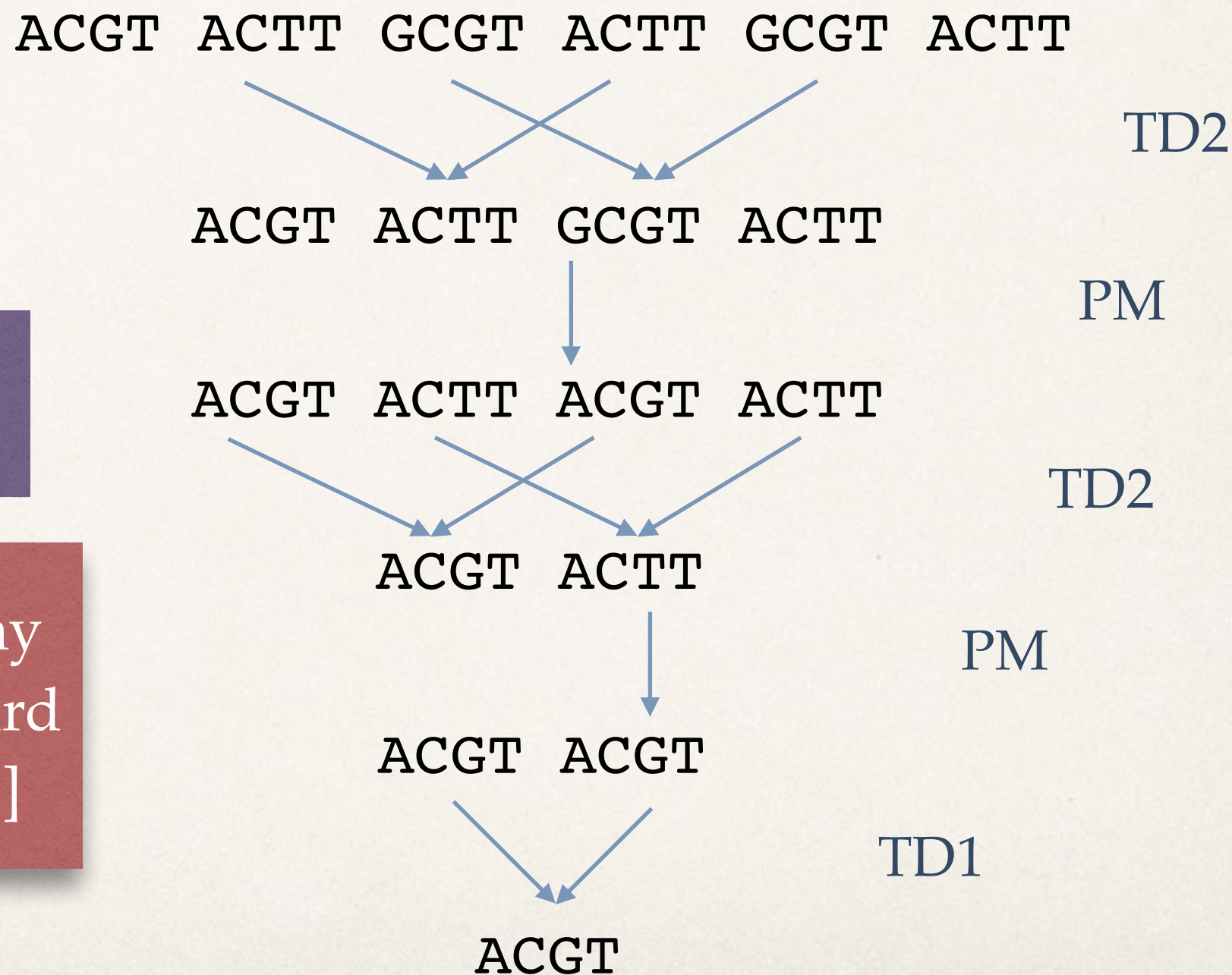
ACGT

TD1, 2 TD2, 2 PM



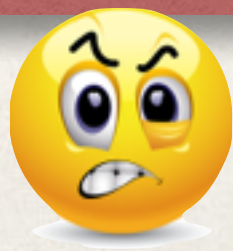


# Finding Duplication History



TD1, 2 TD2, 2 PM

Maximum Parsimony  
Thought to be NP-hard  
[Gascuel et al., 2005]





# Given the final sequence, can we efficiently estimate the parameters?

---

GCT <b>C</b> CGTTACAGGTGGGC <b>A</b> GGGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGGG <b>A</b> GGC <b>G</b>	GCT <b>C</b> CGTTACAGGTGGGC <b>A</b> GGGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGGG <b>A</b> GGC <b>G</b>	GCT <b>C</b> CGTTACAGGTGGGC <b>A</b> GGGGG <b>A</b> GGC <b>G</b>
GCT <b>C</b> CGTTACAGGTGGGC <b>A</b> GGGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGGG <b>G</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>A</b> GGC <b>G</b>	GCT <b>C</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>G</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>A</b> GGC <b>T</b>	GCT <b>C</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>A</b> GGC <b>T</b>
GCT <b>C</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>G</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>G</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>A</b> GGC <b>G</b>
GCT <b>C</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>A</b> GGGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>A</b> GGC <b>G</b>	GCT <b>C</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>A</b> GGC <b>G</b>
GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>A</b> GGC <b>G</b>	GCT <b>G</b> CGTTACAGGTGGGC <b>G</b> GGGGG <b>A</b> GGC <b>G</b>

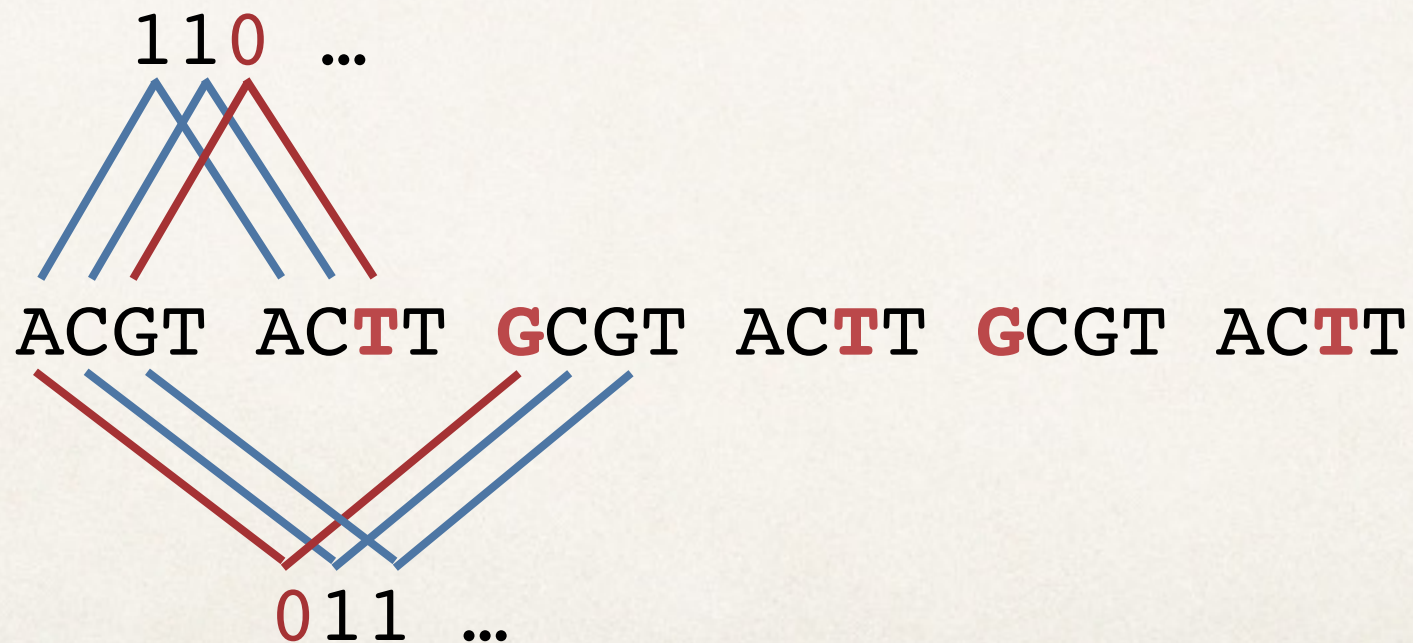


# How to extract information from point mutations?

---

- ❖ Autocorrelation function:

- ★  $r(\delta) = \text{fraction of symbols at distance } \delta \text{ units that are the same}$



$$r(1) = 11/20$$

$$r(2) = 15/16$$



# Stochastic Approximation

---

- ❖ Suppose a discrete random process  $x$  satisfies:

$$x_{n+1} - x_n = \frac{1}{n} (h(x_n) + M_{n+1})$$

for a Lipschitz function  $h$ , and martingale difference  $M$ .

- ❖ Then  $x_n$  converges almost surely to a compact connected internally chain transitive invariant set of the ode

$$\dot{x}_t = h(x_t).$$



# Stochastic Approximation for Autocorrelation

---



# Stochastic Approximation for Autocorrelation

---

- ❖  $r_n$ : autocorr. after  $n$  mutations
- ❖ The stochastic approximation equation for  $r_n$ :

$$\frac{d}{dt}r_t = Ar_t$$

$A$ : a matrix that depends on the parameters:  
 $P(\text{PM}), P(\text{TD1}), P(\text{TD2}), \dots$

- ❖ *As  $n$  increases,  $r_n$  tends to a point in the null space of  $A$*



# Autocorrelation Limit

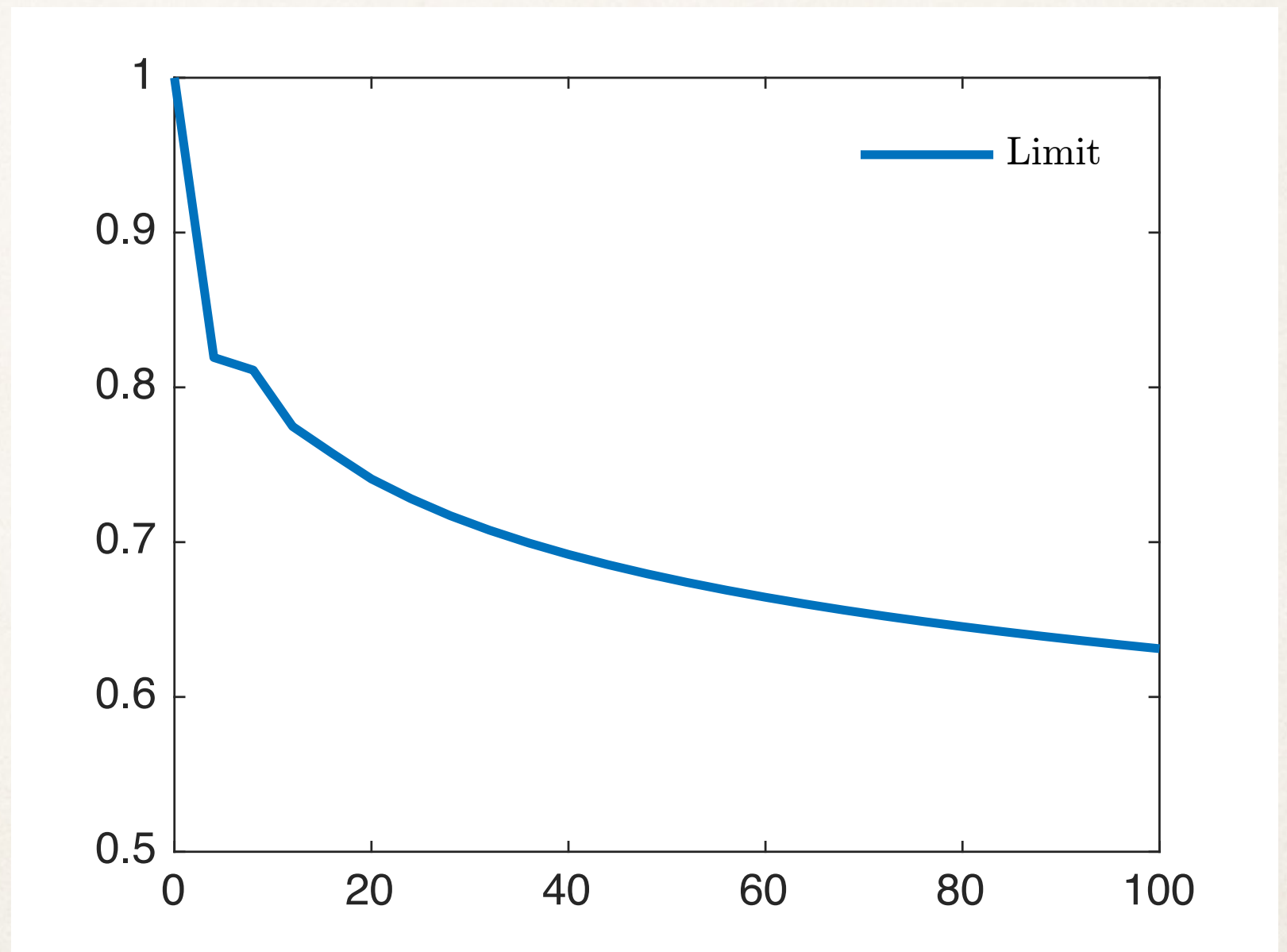
---

$$P(\text{PM}) = 0.250$$

$$P(\text{TD1}) = 0.525$$

$$P(\text{TD2}) = 0.225$$

$r_n(\delta)$



$\delta$



# Autocorrelation Limit

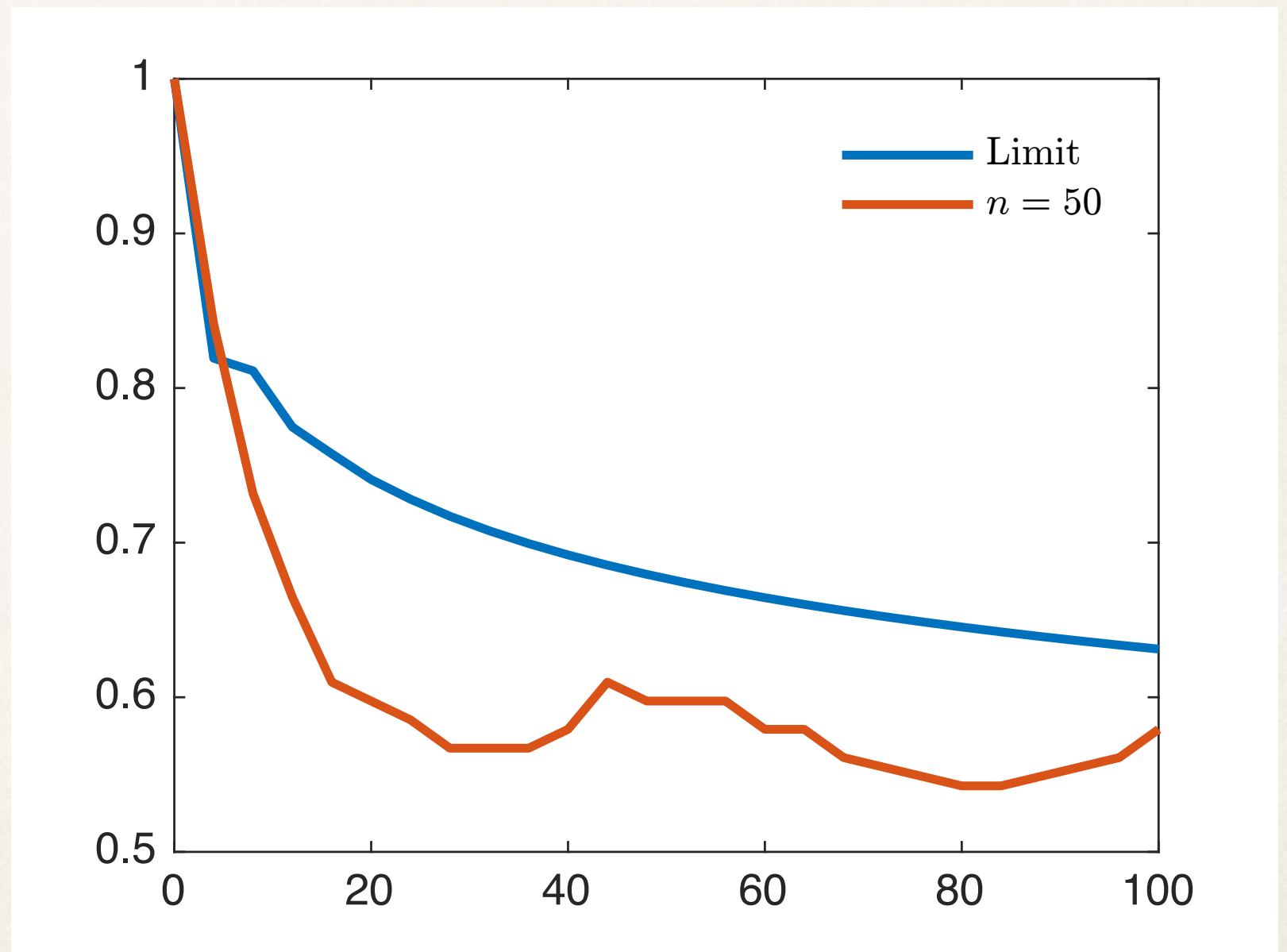
---

$$P(\text{PM}) = 0.250$$

$$P(\text{TD1}) = 0.525$$

$$P(\text{TD2}) = 0.225$$

$r_n(\delta)$



$\delta$



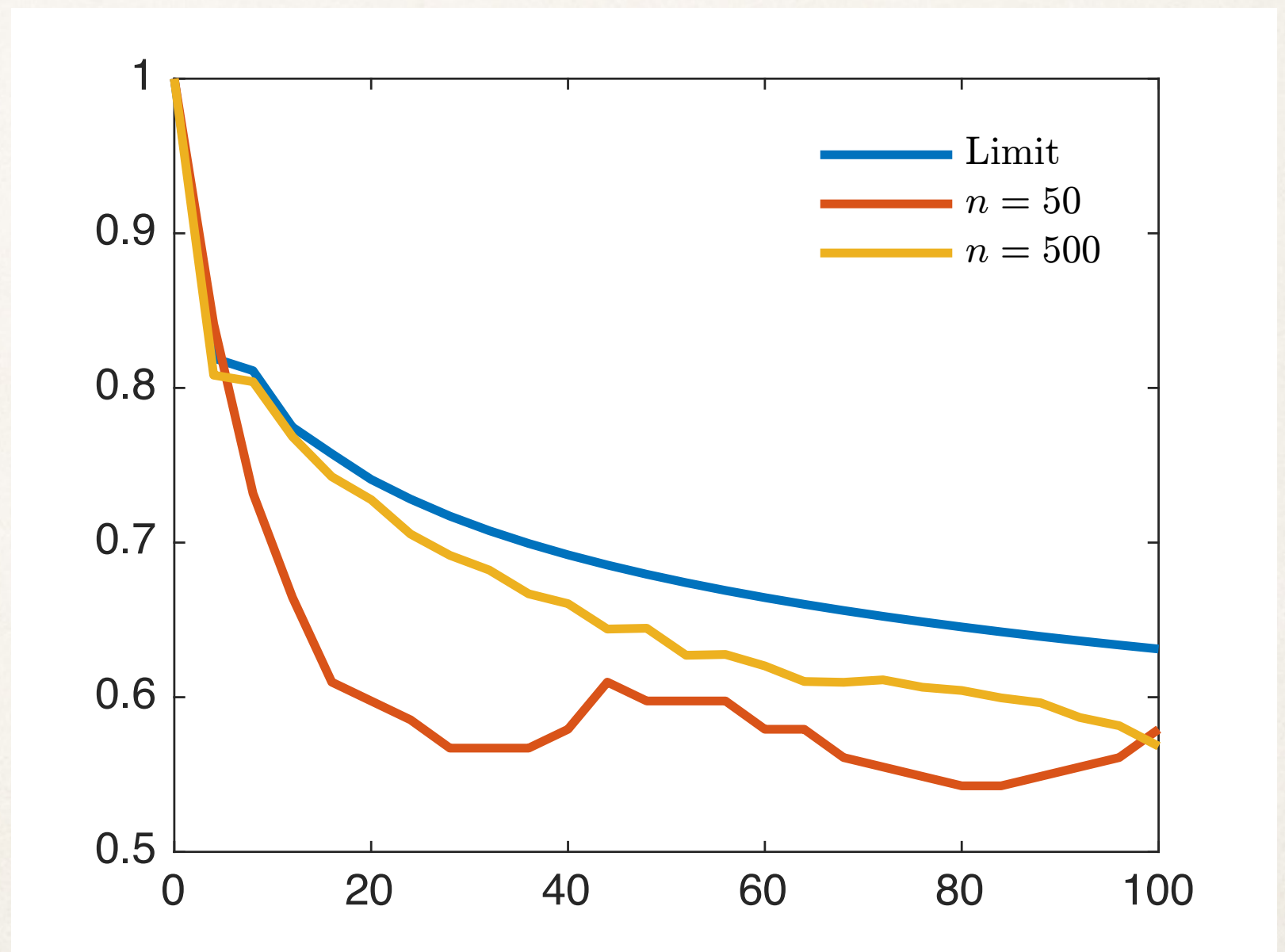
# Autocorrelation Limit

$$P(\text{PM}) = 0.250$$

$$P(\text{TD1}) = 0.525$$

$$P(\text{TD2}) = 0.225$$

$r_n(\delta)$



$\delta$



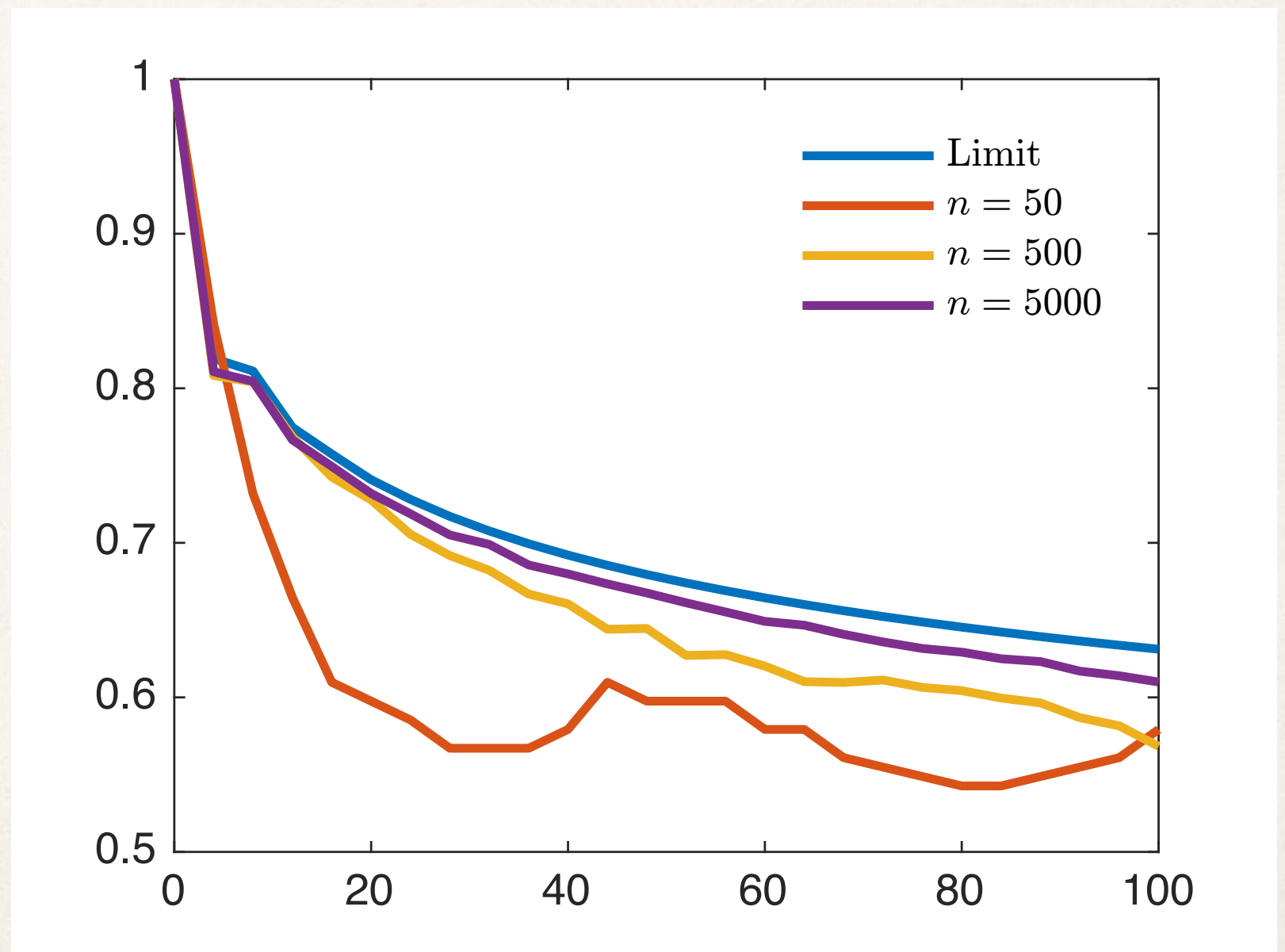
# Autocorrelation Limit

$$P(\text{PM}) = 0.250$$

$$P(\text{TD1}) = 0.525$$

$$P(\text{TD2}) = 0.225$$

$r_n(\delta)$



$\delta$



# Estimation Algorithm

---

$s =$

GCT <b>C</b> CGTTACAGGTGGGCAGGGGAGGCCG	GCT <b>G</b> CGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTCCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTCCGTTACAGGTGGGCAGGGGAGGCCG
GCTCCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTCCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG

1. Calculate autocorrelation  $r$  of  $s$ .
2. Find mutation probs such that the  $l_2$ -norm  $\|Ar\|_2$  is minimized.



# Estimation Algorithm

---

$s =$

GCTCCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTCCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTCCGTTACAGGTGGGCAGGGGAGGCCG
GCTCCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTCCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTCCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTCCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTCCGTTACAGGTGGGCAGGGGAGGCCG
GCTGCGTTACAGGTGGGCAGGGGAGGCCG	GCTGCGTTACAGGTGGGCAGGGGAGGCCG



1. Calculate autocorrelation  $r$  of  $s$ .
2. Find mutation probs such that the  $l_2$ -norm  $\|Ar\|_2$  is minimized.







# Simulation

---



# Simulation

---

- ❖ Start with a short random seed over {A,C,G,T}

TGAATGT

# Simulation

---

- ❖ Start with a short random seed over {A,C,G,T}
- ❖ Choose the parameters  $\mathbf{q} = (PM1, TD1, TD2, TD3)$  randomly

TGAATGT

$$\mathbf{q} = (0.24, 0.33, 0.34, 0.09)$$



# Simulation

---

- ❖ Start with a short random seed over {A,C,G,T}
- ❖ Choose the parameters  $\mathbf{q} = (PM1, TD1, TD2, TD3)$  randomly
- ❖ Apply  $n$  random mutations

TGAATGT

$$\mathbf{q} = (0.24, 0.33, 0.34, 0.09)$$

200 mutations: TGAATGTGCGT...

# Simulation

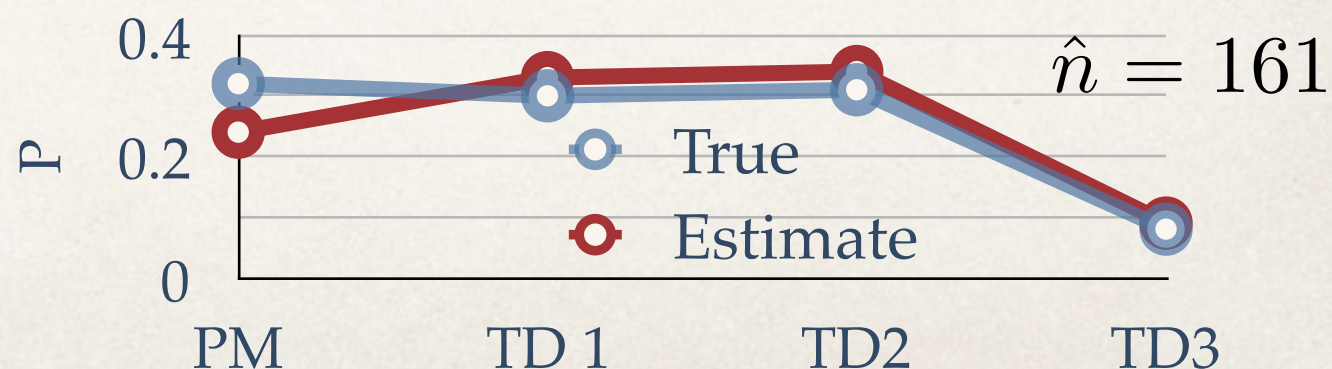
- ❖ Start with a short random seed over  $\{A,C,G,T\}$
- ❖ Choose the parameters  $\mathbf{q} = (PM1, TD1, TD2, TD3)$  randomly
- ❖ Apply  $n$  random mutations
- ❖ Estimate the parameters

$$\min_{\hat{\mathbf{q}}} \|Ar\|_2$$

TGAATGT

$$\mathbf{q} = (0.24, 0.33, 0.34, 0.09)$$

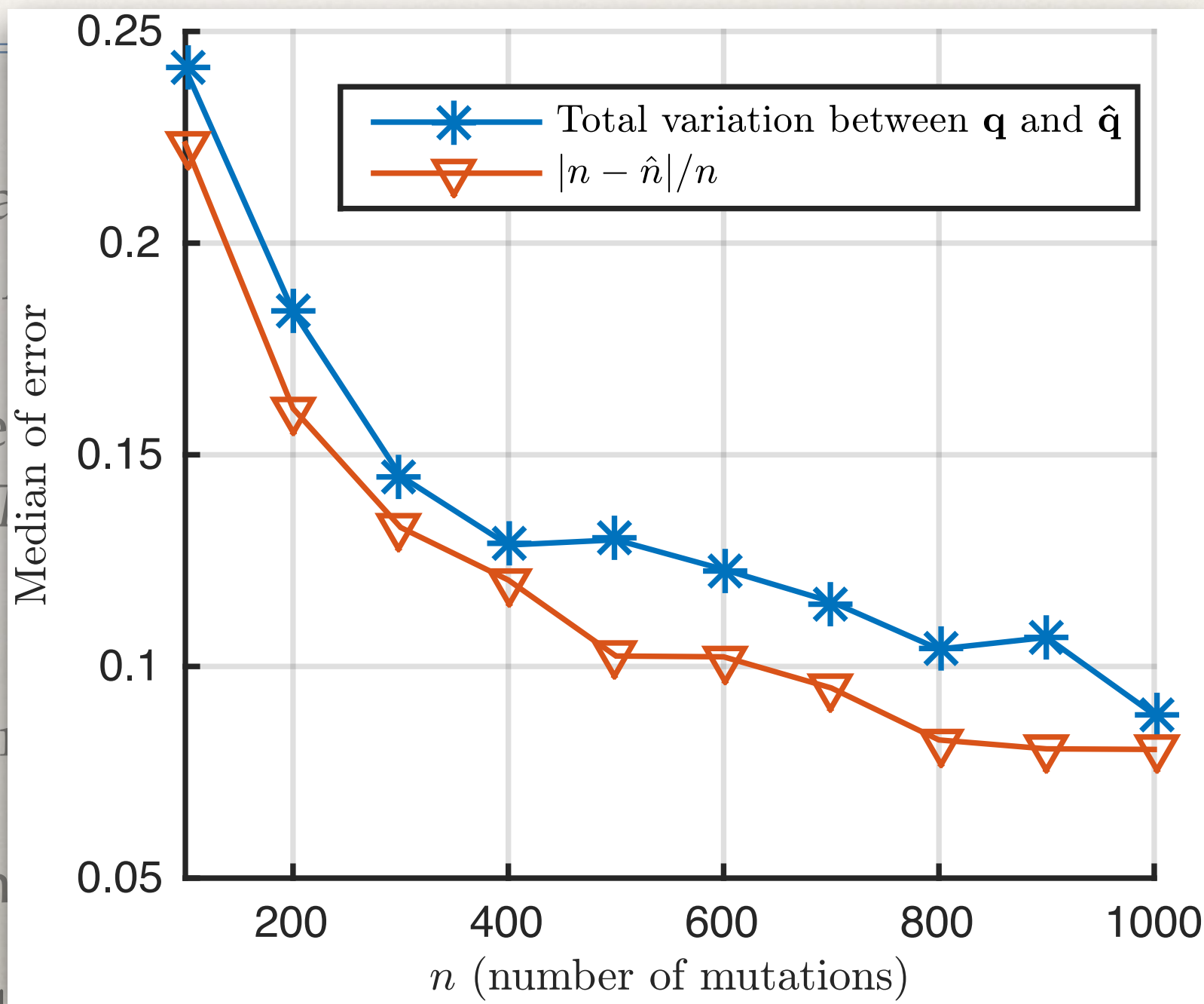
200 mutations: TGAATGTGCGT...





# Simulation

- ❖ Start with a seed over {
- ❖ Choose the  $\mathbf{q} = (PM1, TD1, TD2, TD3)$  randomly
- ❖ Apply  $n$  random mutations
- ❖ Estimate the

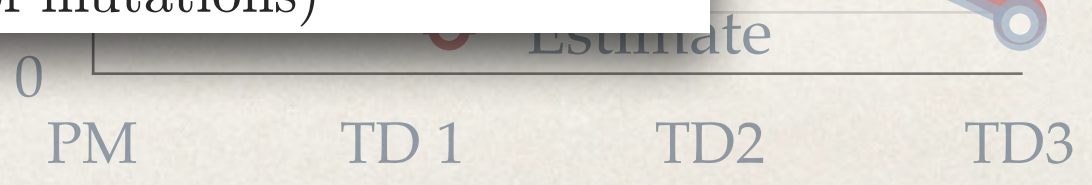


(0.34, 0.09)

ATGTGCGT...

$\hat{n} = 161$

$$\min_{\hat{\mathbf{q}}} \|A\mathbf{r}\|_2$$



# Summary and Next Steps

---



# Summary and Next Steps

---

- ❖ Stochastic estimation algorithm (NP-Hard(?) combinatorial problem).

# Summary and Next Steps

---

- ❖ Stochastic estimation algorithm (NP-Hard(?) combinatorial problem).
- ❖ Point mutation enables estimation of duplication lengths.







