

A photograph of a modern, multi-story glass-walled building at dusk. The building's interior lights are on, and the sky is a mix of purple and blue. In the foreground, there is a circular pond with a concrete edge, reflecting the building and the sky. A grassy area with some trees and a bicycle rack is visible between the pond and the building.

# What's Behind BLAST

*Gene Myers, Director  
MPI for Cell Biology and Genetics  
Dresden, DE*

# Approximate String Search

Given a string  $A$  of length  $n$ , a query  $Q$  of length  $p \ll n$ ,  
an alignment scoring function  $\delta$ , and a threshold  $d$ :

Find all substrings of  $A$ , say  $M$ , s.t.  $\delta(Q, M) \leq d$ ?

$\delta$  here = **Simple Levenstein**

(unit cost mismatch, insert, & delete)

```
...xxxxxxxxaacgt-gcattacxxxxxxxx...  
aatgtggc-ttac
```

A **3-match** (absolute)

A **25%-match** (relative)

# The Story

- **March '88:** The Lister Hill Meeting & Galil's 2 questions

# The Beginning

“Workshop for Algorithms in Molecular Genetics”

March 26-28, 1988

S. Altschul

W. Goad

G. Landau

J. Maisel

T. Smith

M. Zuker

**D. Sankoff**

W. Miller

W. Fitch

T. Hunkapillar

E. Lander

H. Martinez

R. Staden

A. Mukherjee

P. Sellers

G. Myers

Z. Galil

S. Karlin

D. Lipman

C. Sanders

J. Turner

M. Waterman

E. Ukkonen

# Galil's 2 Questions

“Workshop for Algorithms in Molecular Genetics”

March 26-28, 1988

Zvi gave a talk about suffix trees:

**Q1:** Can one get rid of the annoying dependence on alphabet size  $\Sigma$ ?

⇒ Manber & Myers, “Suffix Arrays” 1990

**Q2:** Can one use an index to get faster approximate search?

# The Story

- **March '88:** The Lister Hill Meeting & Galil's 2 questions

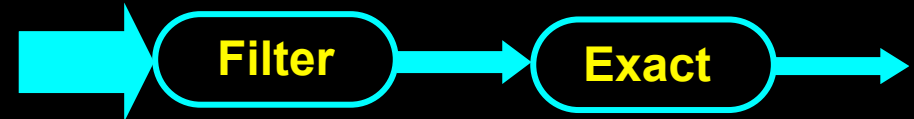
# The Story

- **March '88:** The Lister Hill Meeting & Galil's 2 questions
- **June '88:** Seed & Extend

# APM Filters

A **filter** is an algorithm that eliminates a lot of that which isn't desired.

	100% Sn	< 100% Sn
100% Sp	Exact	
< 100% Sp	Filter	Heuristic



If fast & specific then can improve speed of an exact algorithm.

Approximate match filter ideas:

- Look for exact matches to **k**-mers of the query (in an index)  
( Pearson & Lipman FASTA, Chang & Lawler,  $O(dn/\lg p)$  )
- Instead look for **k**-mers that are a small distance away, e.g. 1 or 2 diff's, from a **k**-mer of the query, i.e. the **neighborhood**

$$\mathcal{N}_d(w) = \{ v : v \text{ and } w \text{ are } \leq d \text{ differences apart} \}$$

$$|\mathcal{N}_d(k)| \approx \binom{k}{d} (2\Sigma)^d$$



# The Power of Neighborhoods

Consider looking for a 9%-match of 40 symbols ( $\Rightarrow \leq 3$  differences or 3-match):



If divide “query” into 4 10-mers then at least one must match exactly:

$\Rightarrow$  Get a hit every  $\Sigma^{10} / 4$  symbols (e.g.  $2.5 \cdot 10^5$  for DNA)



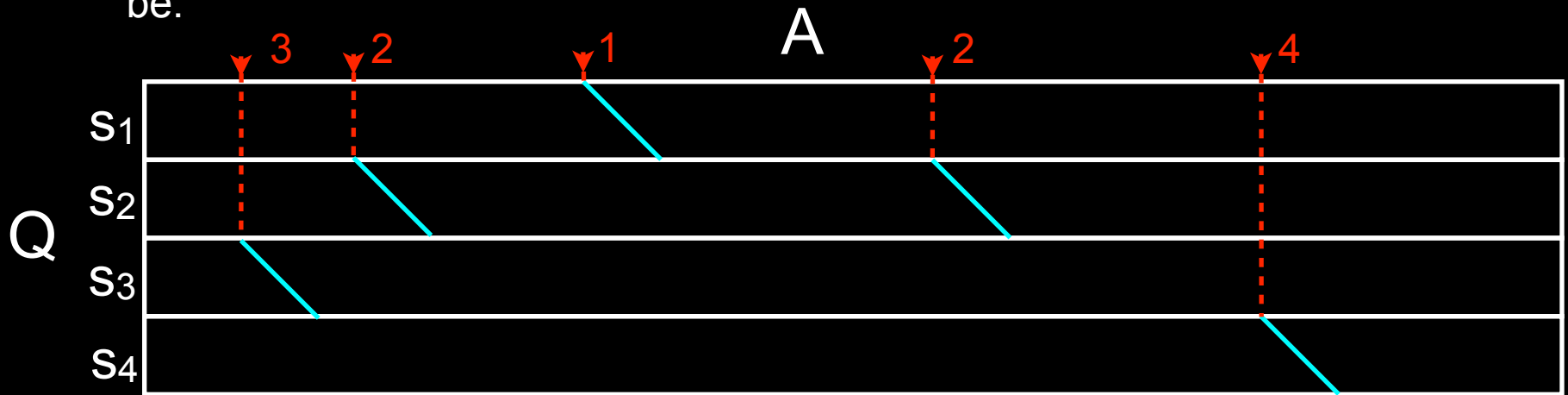
If divide into 2 20-mers then at least one of the  $\mathcal{N}_1$  strings must match exactly:

$\Rightarrow$  Get a hit every  $\Sigma^{20} / 2\mathcal{N}_1(20)$  symbols (e.g.  $10^{12} / 2 \cdot 160 = 3.12 \cdot 10^9$  for DNA)

10,000 times more specific ! (but 80x more lookups)

# “Seed & Extend”

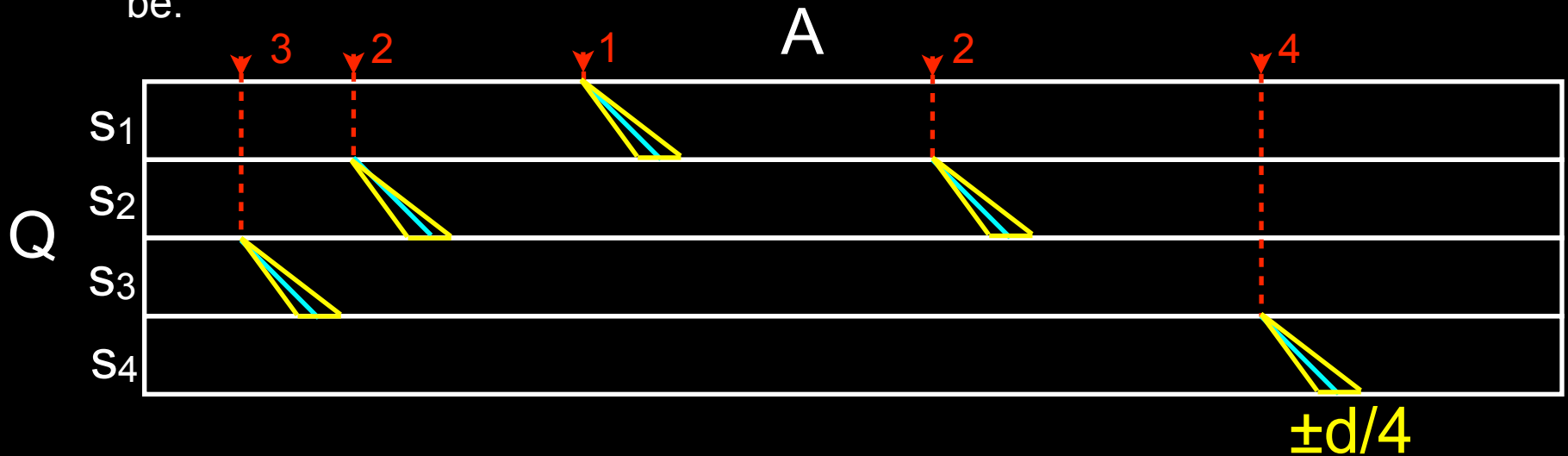
The “seed” matches (either exact or from a neighborhood) are in effect defining areas within the edit graph of  $Q$  vs  $A$  where the alignment of an  $\epsilon$ -match could be:



$$Q = S_1S_2S_3S_4$$

# “Seed & Extend”

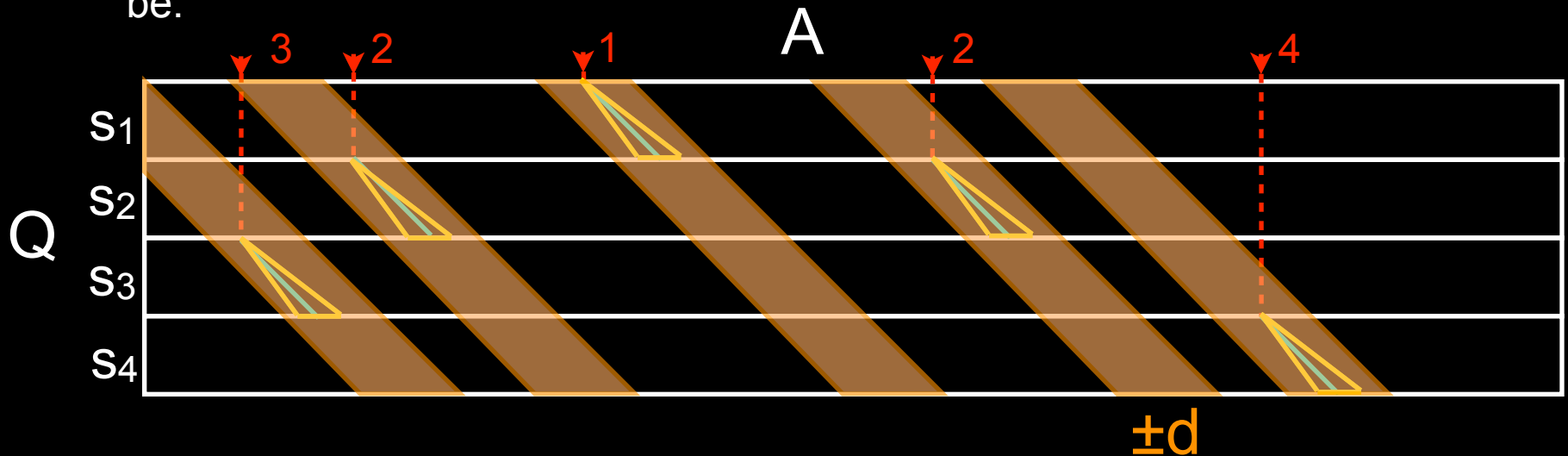
The “seed” matches (either exact or from a neighborhood) are in effect defining areas within the edit graph of  $Q$  vs  $A$  where the alignment of an  $\epsilon$ -match could be:



$$Q = S_1S_2S_3S_4$$

# “Seed & Extend”

The “seed” matches (either exact or from a neighborhood) are in effect defining areas within the edit graph of  $Q$  vs  $A$  where the alignment of an  $\epsilon$ -match could be:



Spend  $O(pdh + pz)$  time where

$h(k)$  = the number of seed “k-hits” vs.  $z(k)$  = neighborhood size “k-words”

Both  $z$  and  $h$  are functions of  $k$  and the optimal  $k$  is “slightly bigger” than  $\log_{\Sigma} n$

# The Story

- **March '88:** The Lister Hill Meeting & Galil's 2 questions
- **June '88:** Seed & Extend

# The Story

- **March '88:** The Lister Hill Meeting & Galil's 2 questions
- **June '88:** Seed & Extend
- **May '89:** The TRW Chip & The Cigarette Break
- **Fall '89:** Blast is Born

## Blast = Seed & Extend

Seeds are neighborhoods of all k-mers of query  
under weighted Levenstein (e.g. PAM120)

Find seeds with a deterministic finite automaton  
accepting all neighborhood words ( $\Rightarrow O(n)$ )

Extend is just weighted Hamming  
but stop when score drops too much

A heuristic

“blast” was inspired by  
“slam” = sublinear approximate match

# The Story

- **March '88:** The Lister Hill Meeting & Galil's 2 questions
- **June '88:** Seed & Extend
- **May '89:** The TRW Chip & The Cigarette Break
- **Fall '89:** Blast is Born



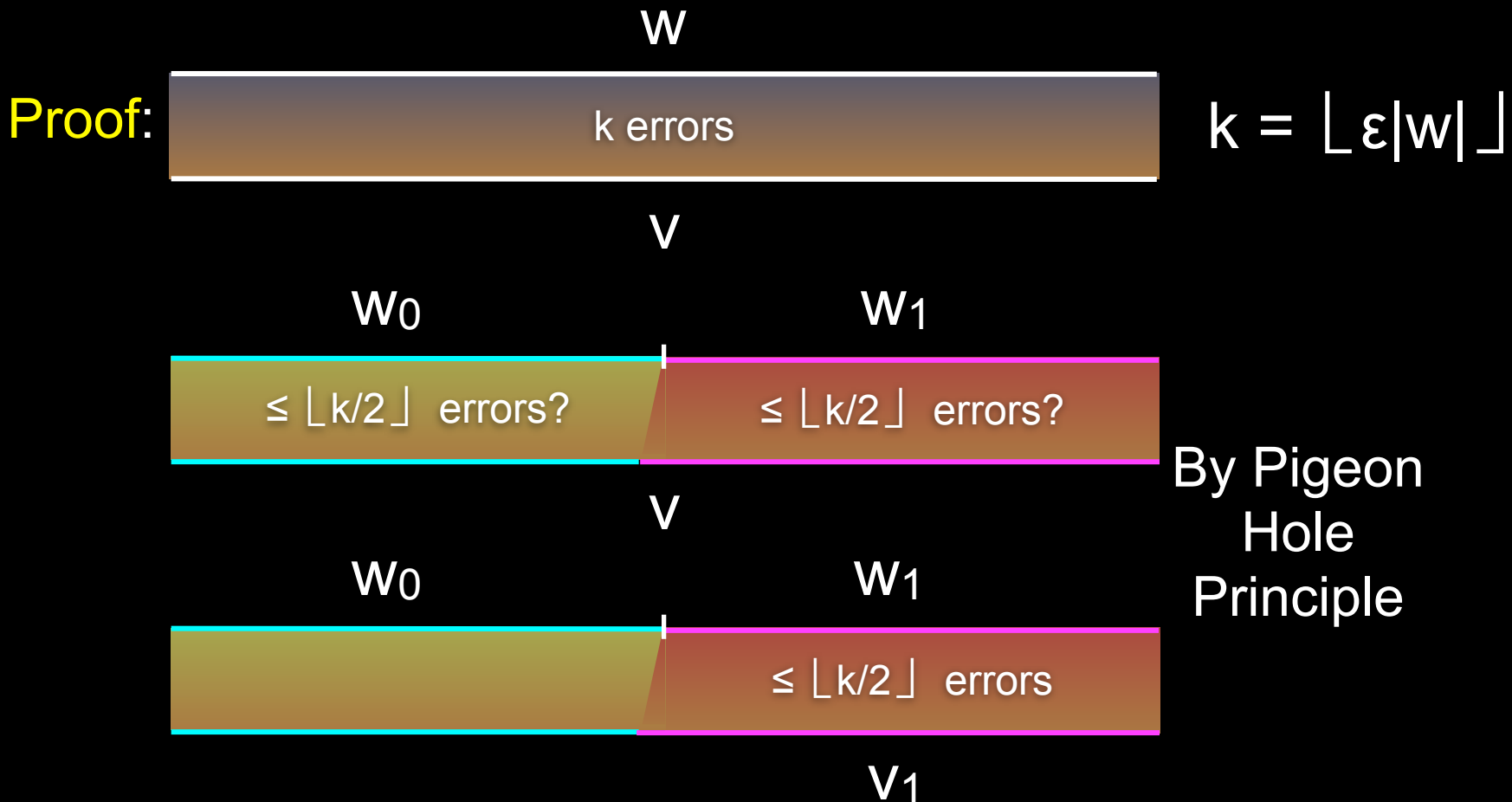
# The Story

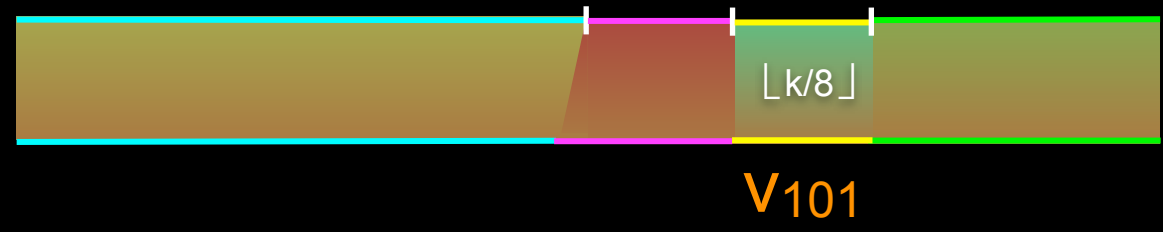
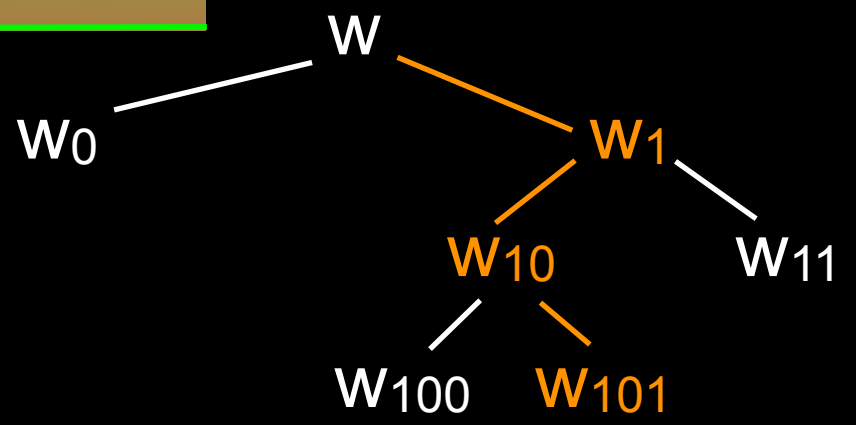
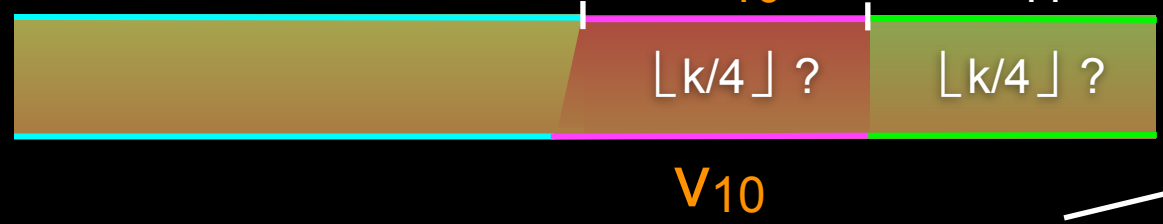
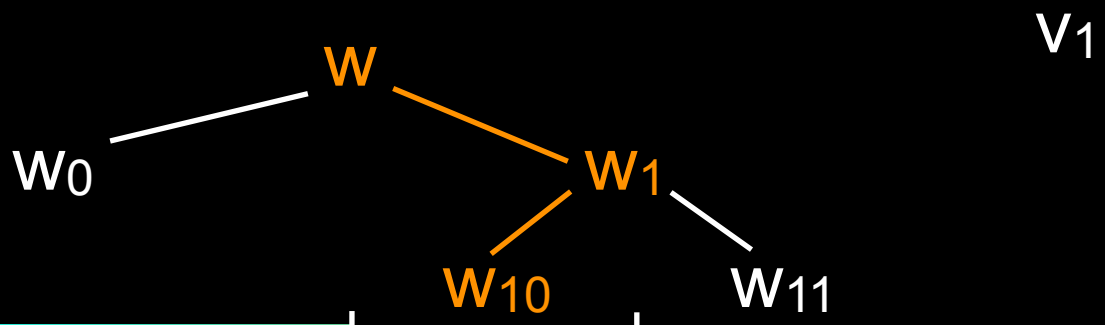
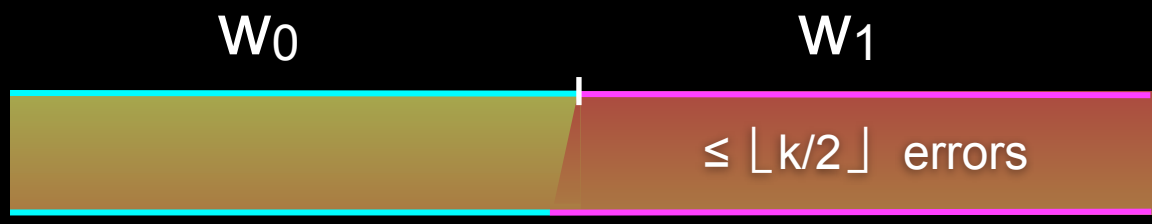
- **March '88:** The Lister Hill Meeting & Galil's 2 questions
- **June '88:** Seed & Extend
- **May '89:** The TRW Chip & The Cigarette Break
- **Fall '89:** Blast is Born
- **Fall '89:** The Splitting Lemma

# The Splitting Lemma

**Lemma:** If  $w$   $\epsilon$ -matches  $v$  then either

- (a)  $w_0$  has an  $\epsilon$ -match to a prefix (call it  $v_0$ ) of  $v$ , or
- (b)  $w_1$  has an  $\epsilon$ -match to a suffix (call it  $v_1$ ) of  $v$ .

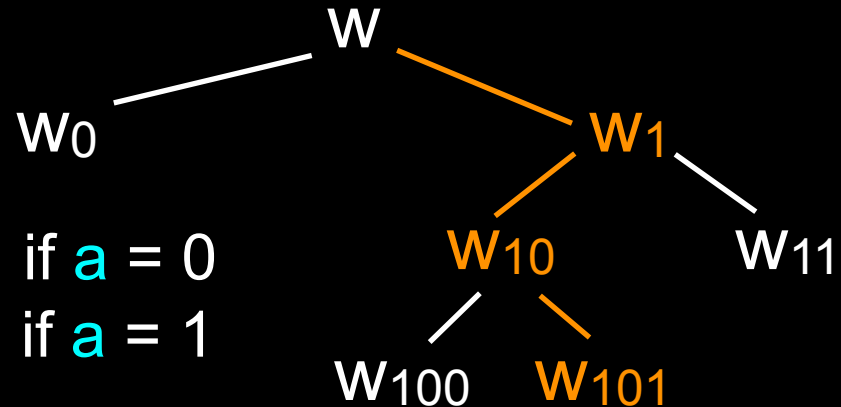




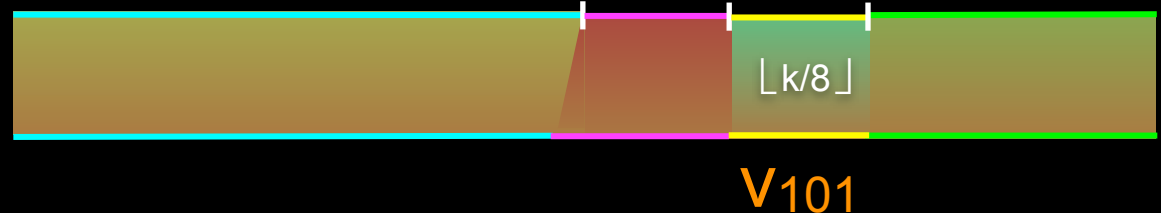
# The Splitting Lemma

Let  $w_\epsilon = w$

$$w_{\beta a} = \begin{cases} w_\beta[1..|w_\beta|/2] & \text{if } a = 0 \\ w_\beta[|w_\beta|/2+1..|w_\beta|] & \text{if } a = 1 \end{cases}$$



e.g.  $\alpha = 101\dots$



- Lemma:** If  $w$   $\epsilon$ -matches  $v$  then  $\exists \alpha$  s.t.  $\forall$  prefixes  $\beta$  of  $\alpha$ ,
- (1)  $w_\beta$  has an  $\epsilon$ -match to a substring (call it  $v_\beta$ ) of  $v$ , and
  - (2)  $v_{\beta 0}$  is a prefix of  $v_\beta$  (if  $\beta 0$  is a prefix of  $\alpha$ ), and
  - (3)  $v_{\beta 1}$  is a suffix of  $v_\beta$  (if  $\beta 1$  is a prefix of  $\alpha$ ).

# The Story

- **March '88:** The Lister Hill Meeting & Galil's 2 questions
- **June '88:** Seed & Extend
- **May '89:** The TRW Chip & The Cigarette Break
- **Fall '89:** Blast is Born
- **Fall '89:** The Splitting Lemma

# The Story

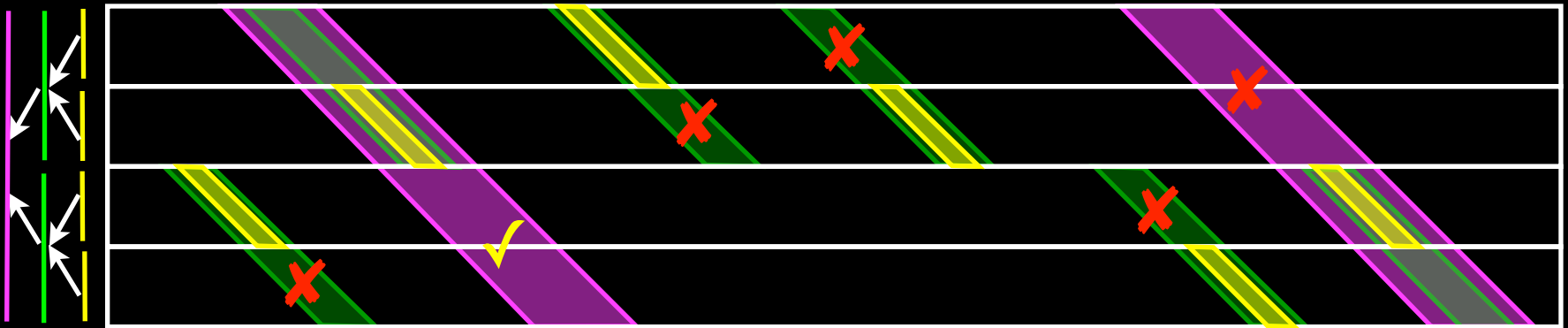
- **March '88:** The Lister Hill Meeting & Galil's 2 questions
- **June '88:** Seed & Extend
- **May '89:** The TRW Chip & The Cigarette Break
- **Fall '89:** Blast is Born
- **Fall '89:** The Splitting Lemma
- **Fall '89:** Seed & Extend by Doubling

# Doubling Extension

Use  $\log_{\Sigma} n$  as the seed size !

**Lemma:** Any  $\epsilon$ -match of  $Q$  has an  $\epsilon$ -match to at least one seed segment of size  $\log_{\Sigma} n$

Use the splitting lemma to split  $Q$  to seeds of size  $\log_{\Sigma} n$ , and instead of extending all at once, extend by doubling using the splitting lemma.



Time for each extension telescopes hyper-geometrically and so is dominated by the first term:

$$O(P/\log_{\Sigma} n \cdot h \cdot \log_{\Sigma} n \cdot \epsilon \log_{\Sigma} n) = O(dh \log_{\Sigma} n)$$

# The Story

- **March '88:** The Lister Hill Meeting & Galil's 2 questions
- **June '88:** Seed & Extend
- **May '89:** The TRW Chip & The Cigarette Break
- **Fall '89:** Blast is Born
- **Fall '89:** The Splitting Lemma
- **Fall '89:** Seed & Extend by Doubling



# The Story

- **March '88:** The Lister Hill Meeting & Galil's 2 questions
- **June '88:** Seed & Extend
- **May '89:** The TRW Chip & The Cigarette Break
- **Fall '89:** Blast is Born
- **Fall '89:** The Splitting Lemma
- **Fall '89:** Seed & Extend by Doubling
- **Spr '90:** Generating Condensed Neighborhoods

# Generating (Condensed) Neighborhoods

$\overline{\mathcal{N}}_d(w) = \{ v : v \text{ and } w \text{ are } \leq d \text{ differences apart and } v \text{ is not a proper prefix of another word in } N_d(w) \}$

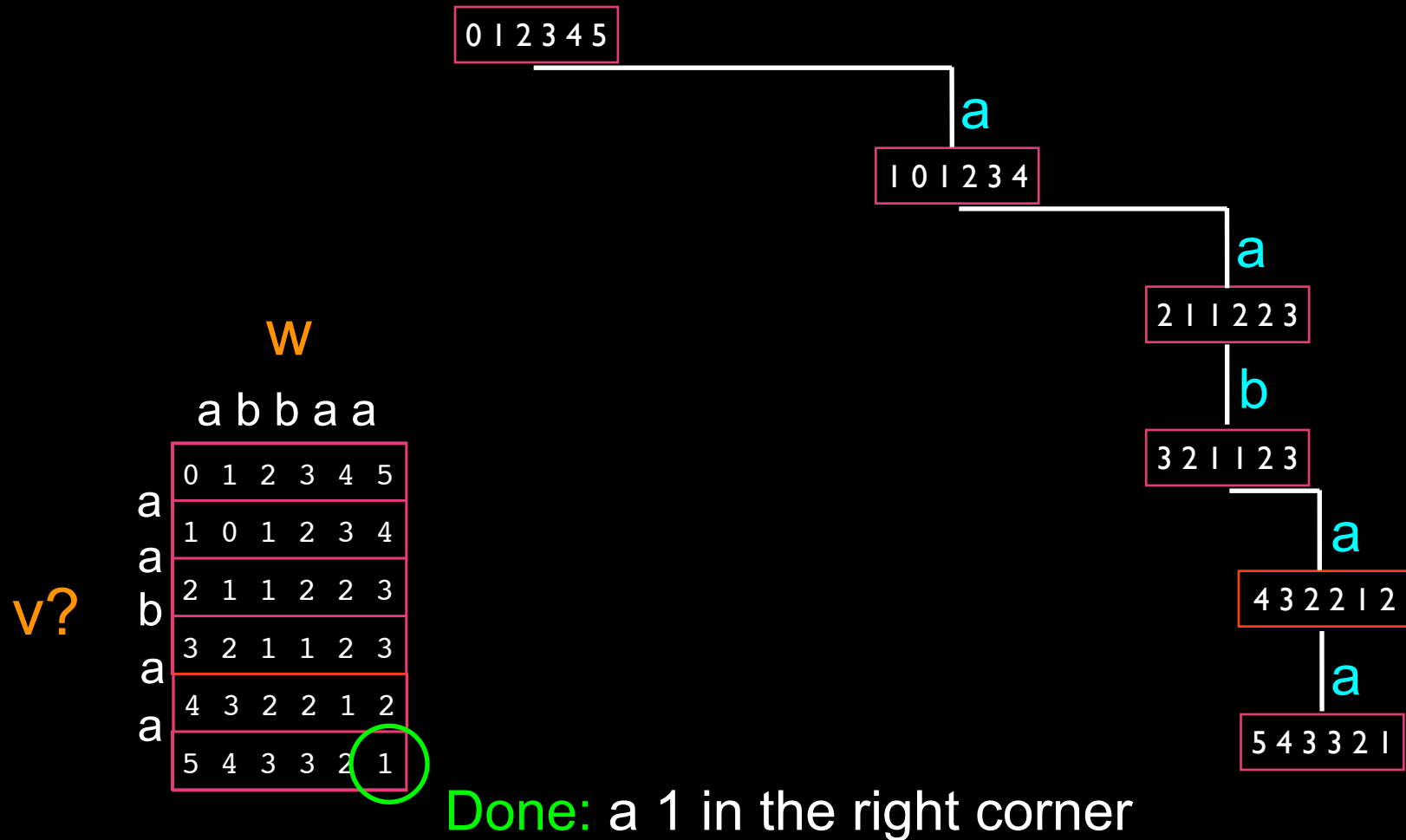
$\overline{\mathcal{N}}_1(\text{abbaa}) = \{ \text{aabaa}, \text{aabbaa}, \text{abaa}, \text{abaaa}, \text{ababaa}, \text{abba}, \text{abbbaa}, \text{abbab}, \text{abbbaaa}, \text{abbaba}, \text{abbba}, \text{abbbaa}, \text{babbaa}, \text{bbaa}, \text{bbbaa} \}$

It suffices to find the words in the condensed neighborhood.

But how do you do that efficiently, including finding them in the index? ...

... Compute rows of dynamic programming matrix as one traverses the trie of all strings over  $\Sigma$

# Condensed Neighborhoods

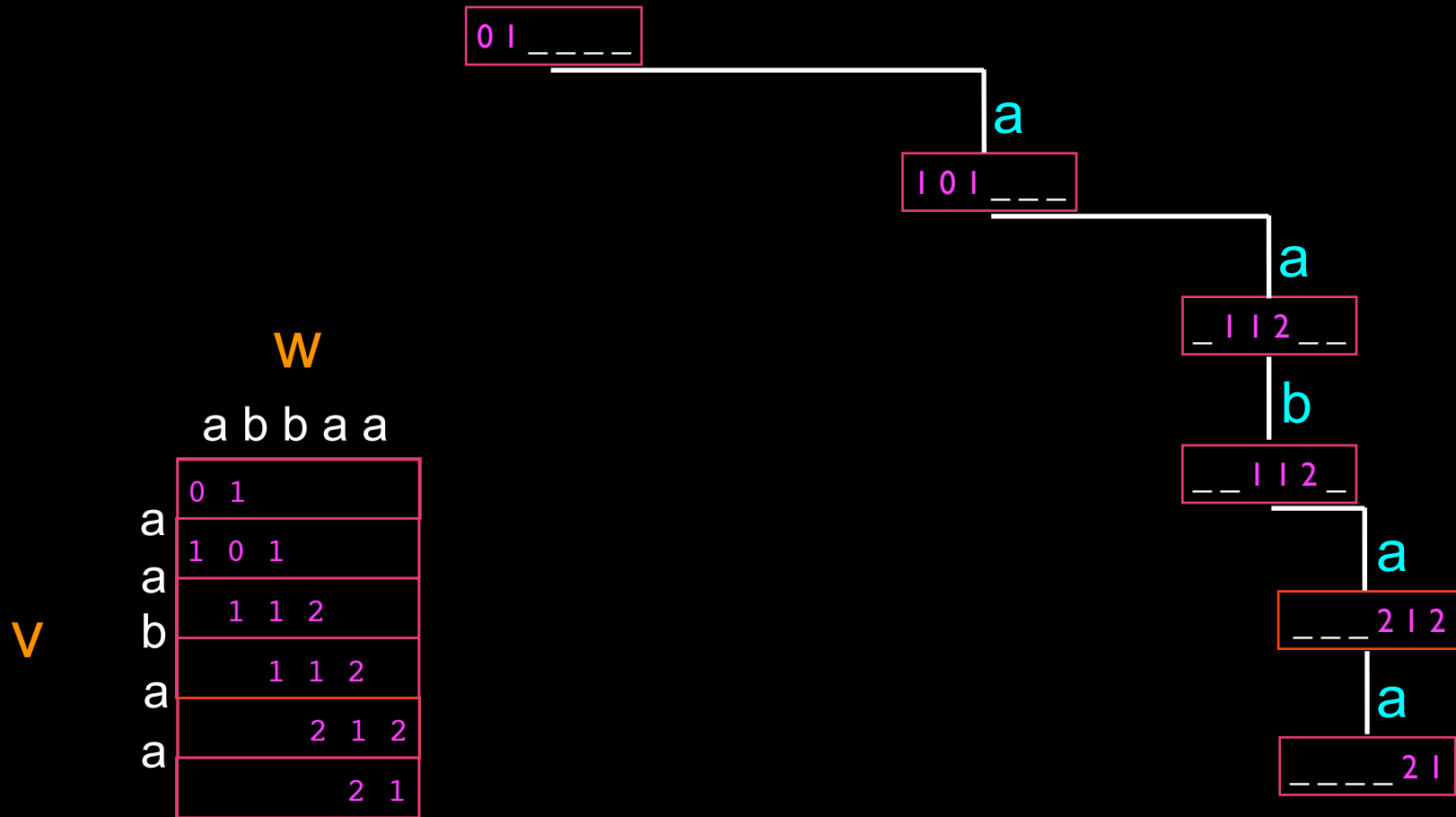


# Condensed Neighborhoods



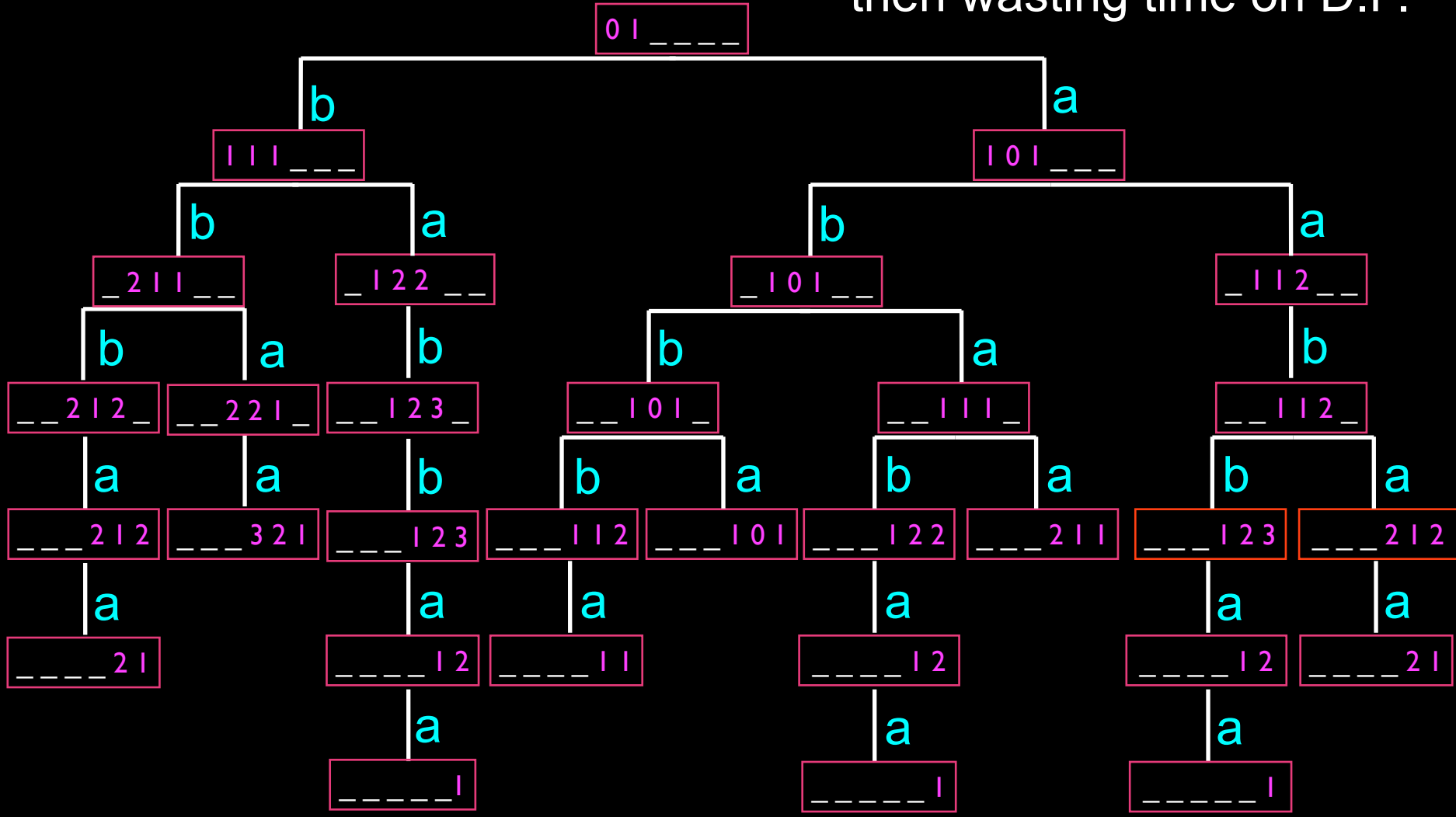
Only need  $\pm 1$  band !

# Condensed Neighborhoods



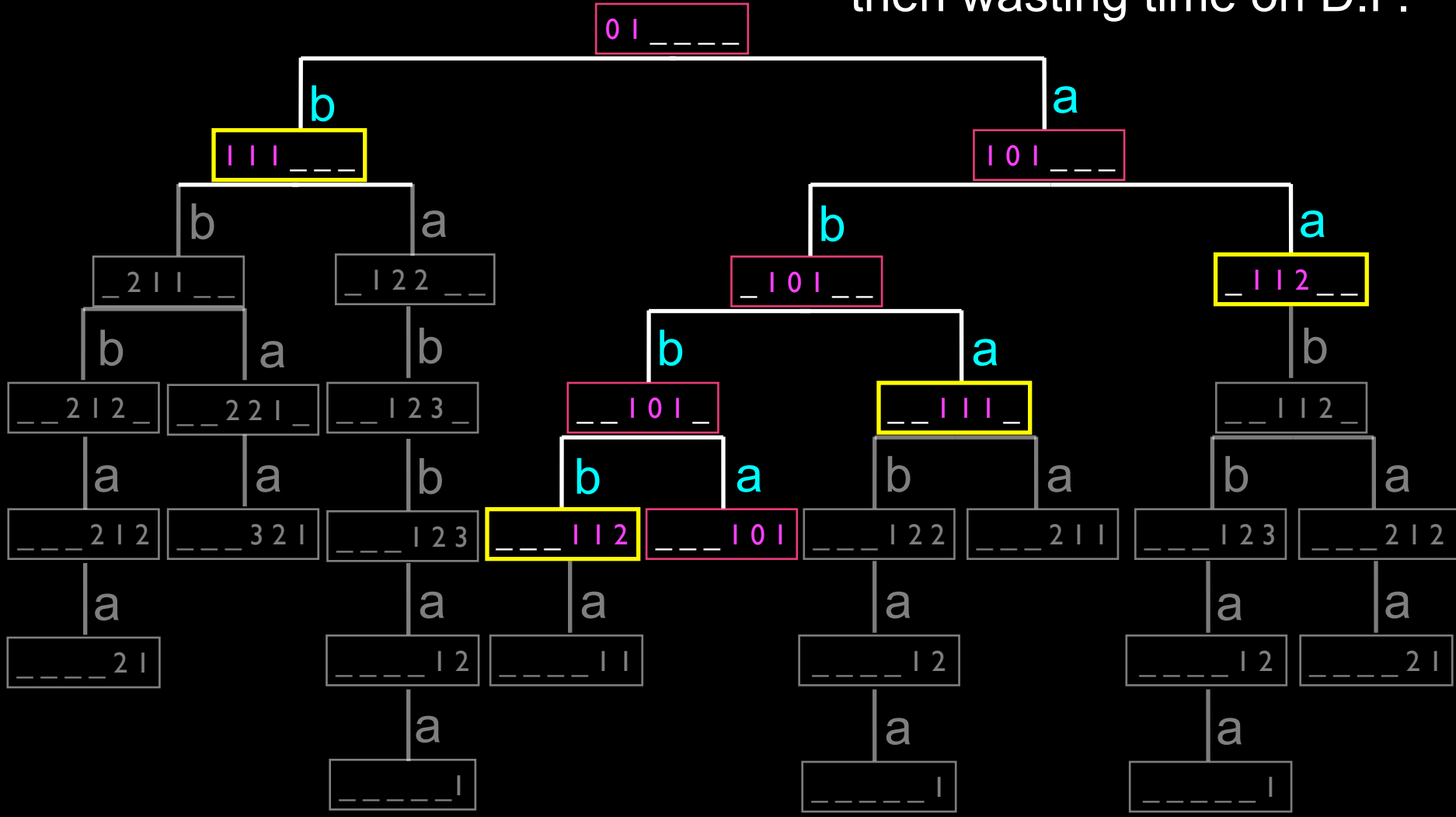
# Condensed Neighborhoods

If all entries are  $\geq d$   
then wasting time on D.P.

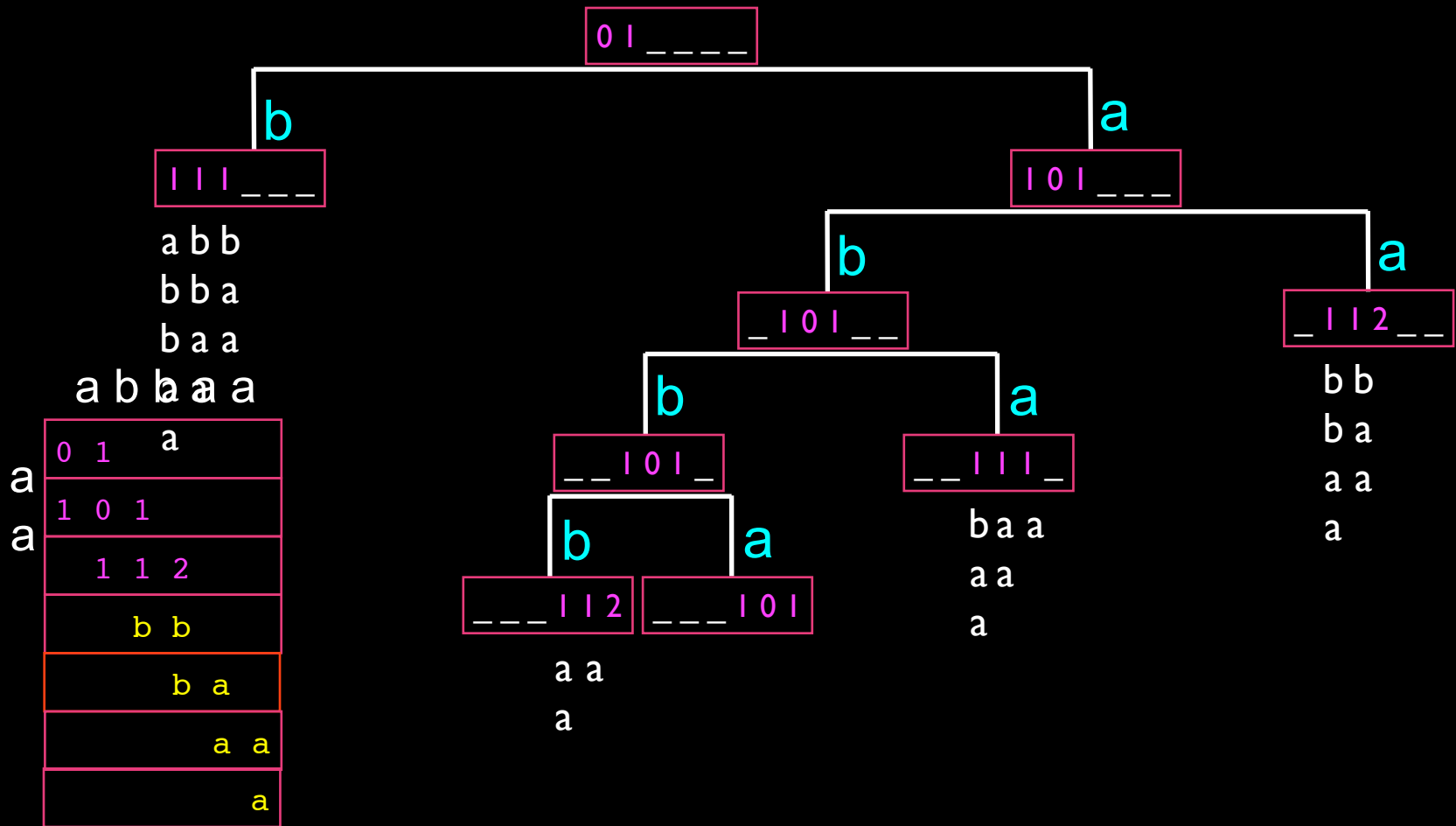


# Condensed Neighborhoods

If all entries are  $\geq d$   
then wasting time on D.P.

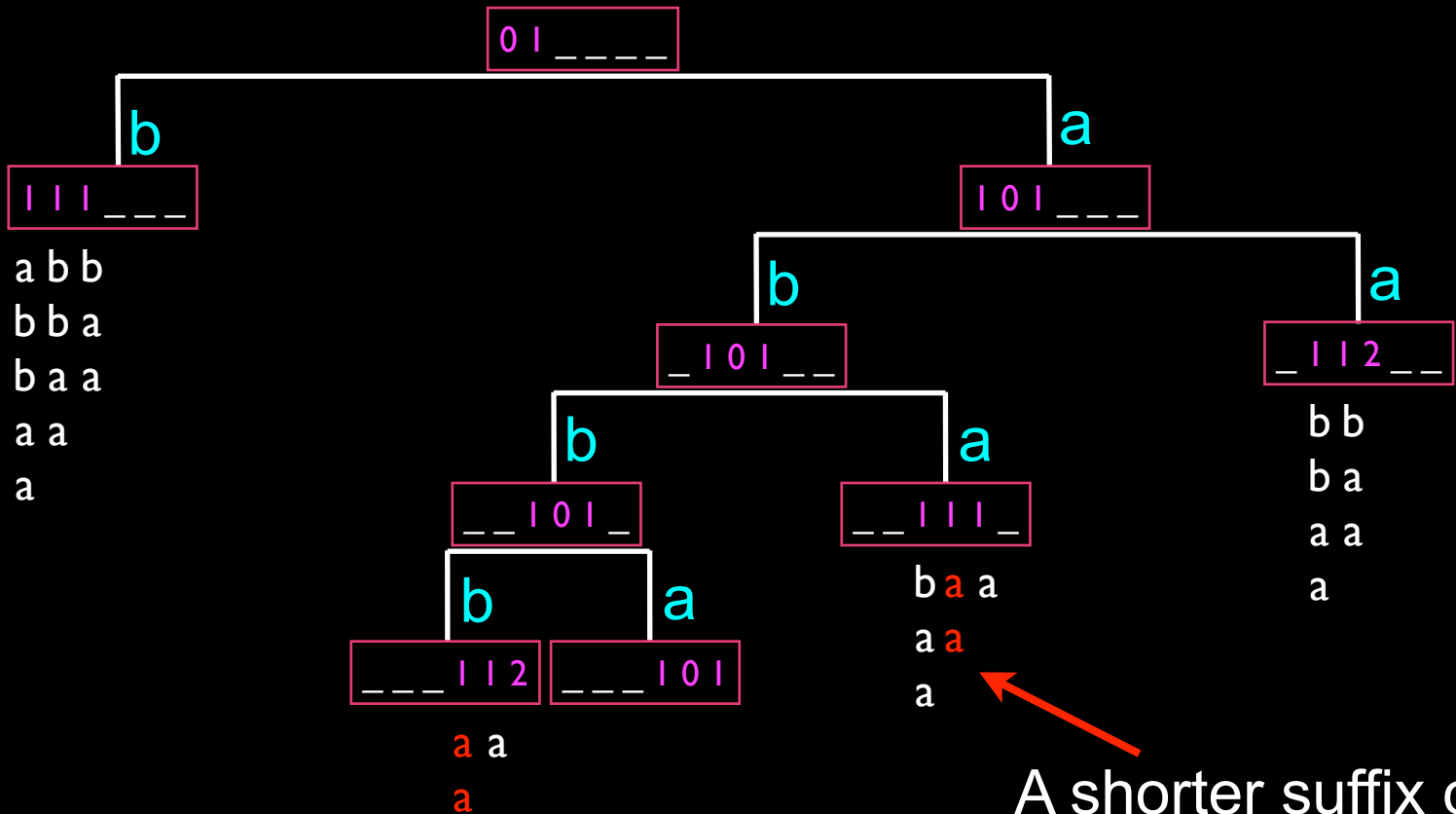


# Condensed Neighborhoods



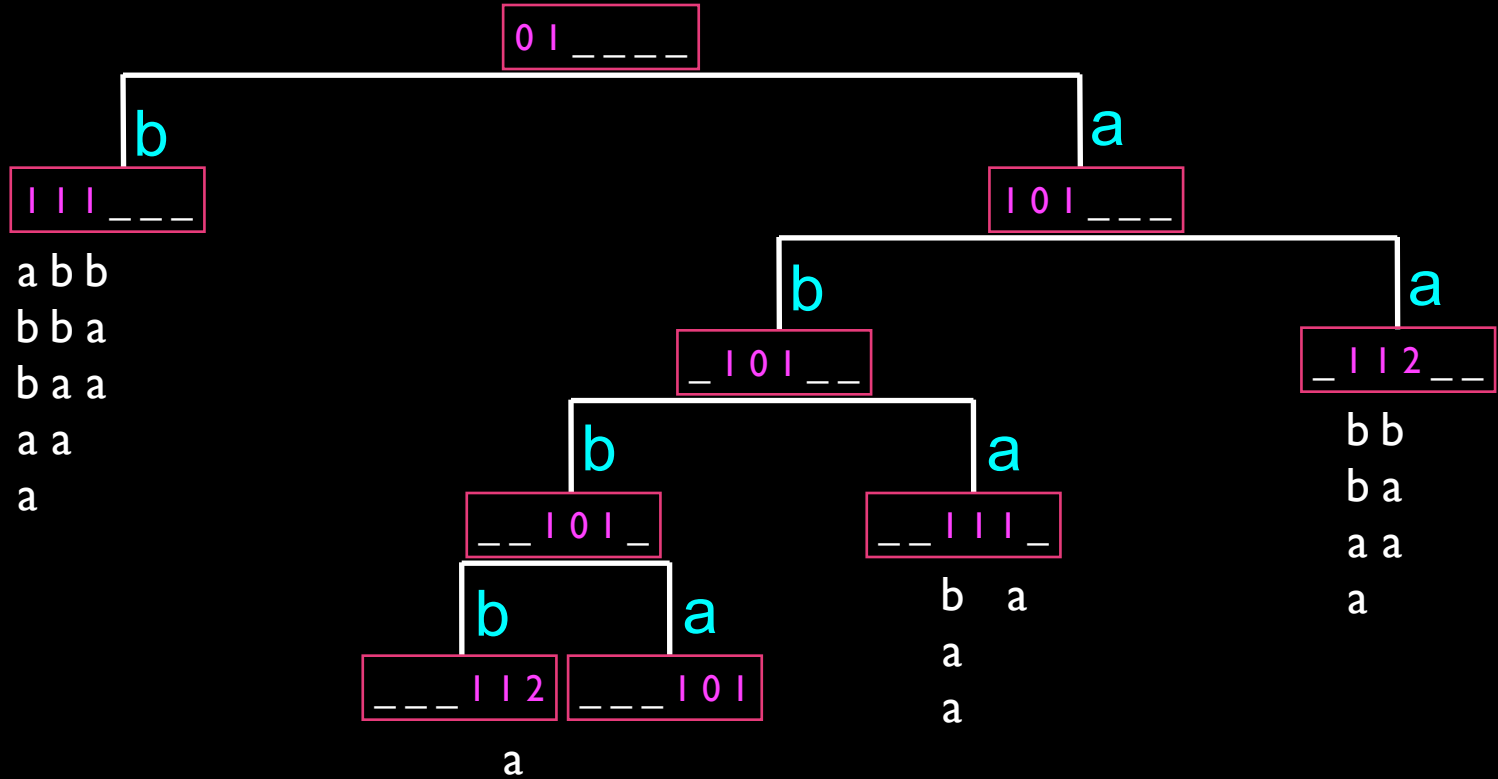


# Condensed Neighborhoods



Use “KMP” on reverse of  $w$  to efficiently discover these.

# Condensed Neighborhoods



**Lemma:** Neighborhoods and their hits in  $A$  can be generated in  $O(zd+h)$  time where  $z = |\mathcal{N}_d(w)|$

# The Story

- **March '88:** The Lister Hill Meeting & Galil's 2 questions
- **June '88:** Seed & Extend
- **May '89:** The TRW Chip & The Cigarette Break
- **Fall '89:** Blast is Born
- **Fall '89:** The Splitting Lemma
- **Fall '89:** Seed & Extend by Doubling
- **Spr '90:** Generating Condensed Neighborhoods

# The Story

- **March '88:** The Lister Hill Meeting & Galil's 2 questions
- **June '88:** Seed & Extend
- **May '89:** The TRW Chip & The Cigarette Break
- **Fall '89:** Blast is Born
- **Fall '89:** The Splitting Lemma
- **Fall '89:** Seed & Extend by Doubling
- **Spr '90:** Generating Condensed Neighborhoods
- **Fall '90:** Finale: Complexity

# Complexity

How big is  $\overline{\mathcal{N}}_d(k)$ ?

Developed recurrence for non-redundant edit scripts:

- (a) DI = S
- (b) DS = SD
- (c) IS = SI
- (d) ID =  $\Phi$

Lemma:

$$\begin{aligned} \mathbf{S}(k,d) = & \mathbf{S}(k-1,d) + (\Sigma-1)\mathbf{S}(k-1,d-1) + (\Sigma-1) \sum_{j=0}^{d-1} \Sigma^j \mathbf{S}(k-1,d-1) \\ & + (\Sigma-1)^2 \sum_{j=0}^{d-2} \Sigma^j \mathbf{S}(k-2,d-2-j) + \sum_{j=0}^{d-1} \mathbf{S}(k-2-j,d-1-j) \end{aligned}$$

$$\overline{\mathcal{N}}_d(k) \leq \mathbf{S}(k,d) + \sum_{j=1}^d \Sigma^j \mathbf{S}(k-1,d-j)$$

# Complexity

So how big is it?

Lemma:

$$\overline{\mathcal{N}}_{\varepsilon}(k) \leq 1.708 \alpha(\varepsilon)^k$$

$$\text{where } \alpha(\varepsilon) = \Sigma^{\text{pow}(\varepsilon)}$$

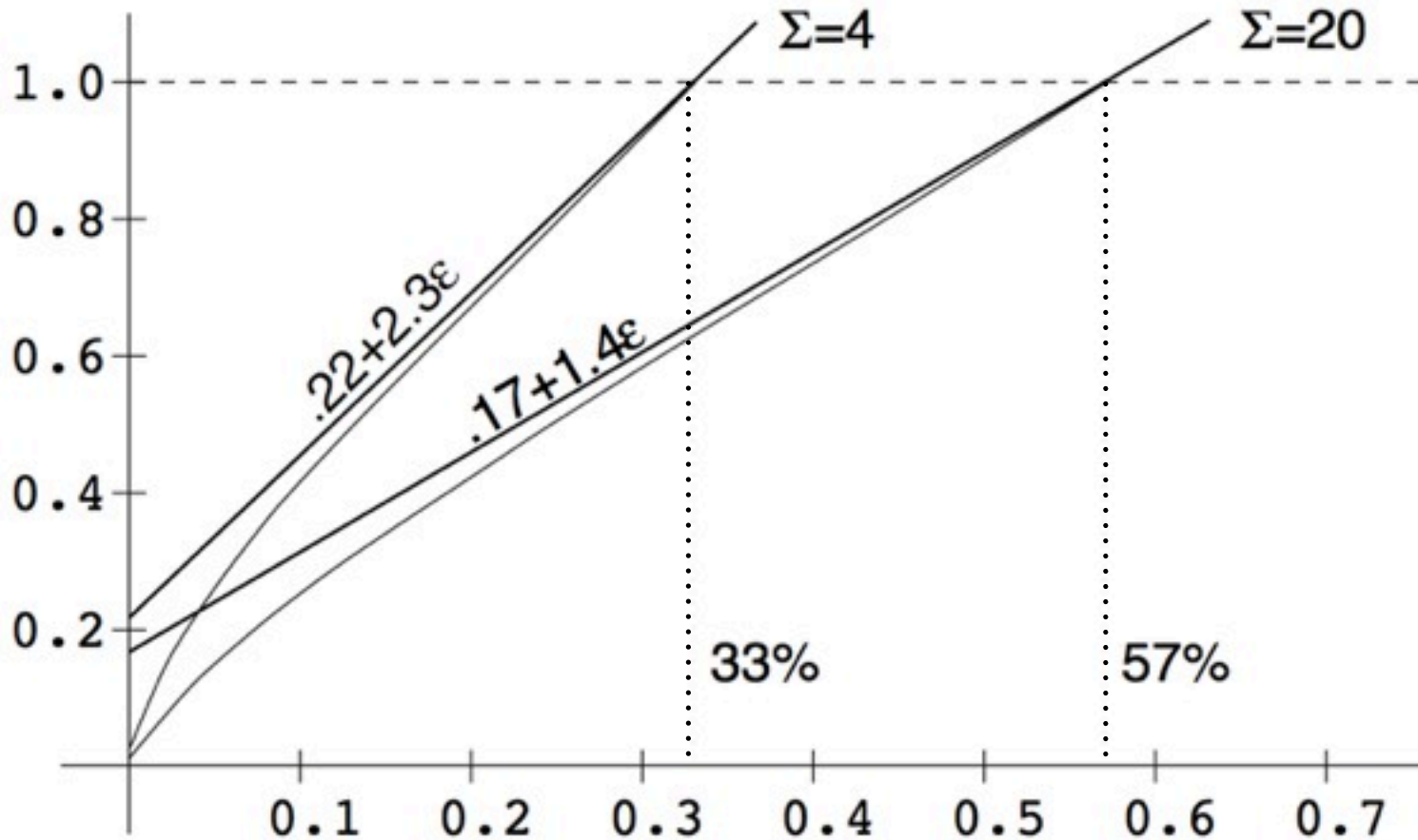
$$\text{and } \text{pow}(\varepsilon) = \log_{\Sigma} \left( \frac{c(\varepsilon)+1}{c(\varepsilon)-1} \right) + \varepsilon \log_{\Sigma} c(\varepsilon) + \varepsilon$$

$$\text{and } c(\varepsilon) = \varepsilon^{-1} + (1 + \varepsilon^{-2})^{.5}$$

$$\text{Also } Pr(w \text{ in } \overline{\mathcal{N}}_{\varepsilon}(k)) = O(1 / \beta(\varepsilon)^k)$$

$$\text{where } \beta(\varepsilon) = \Sigma^{1-\text{pow}(\varepsilon)}$$

$pow(\varepsilon)$



**Figure 4:** Plot of  $pow(\varepsilon)$  and Sample Bounding Lines

# Complexity

So how big is it?

Lemma:

$$\overline{\mathcal{N}}_{\epsilon}(k) = O(\alpha^k)$$

$$\text{where } \alpha = \Sigma^{\text{pow}(\epsilon)}$$

← Starts at 1 ( $\epsilon=0$ ) and grows  
“Flex factor”

$$\Pr(w \text{ in } \overline{\mathcal{N}}_{\epsilon}(k)) = O(1 / \beta^k)$$

$$\text{where } \beta = \Sigma^{1-\text{pow}(\epsilon)}$$

← Starts at  $\Sigma$  ( $\epsilon=0$ ) and shrinks  
“Effective alphabet size”

And when  $k = \log_{\Sigma} n$ ?

$$\overline{\mathcal{N}}_{\epsilon}(k) = O(n^{\text{pow}(\epsilon)}) \quad \text{and} \quad \Pr(w \text{ in } \overline{\mathcal{N}}_{\epsilon}(k)) = O(n^{\text{pow}(\epsilon)-1})$$

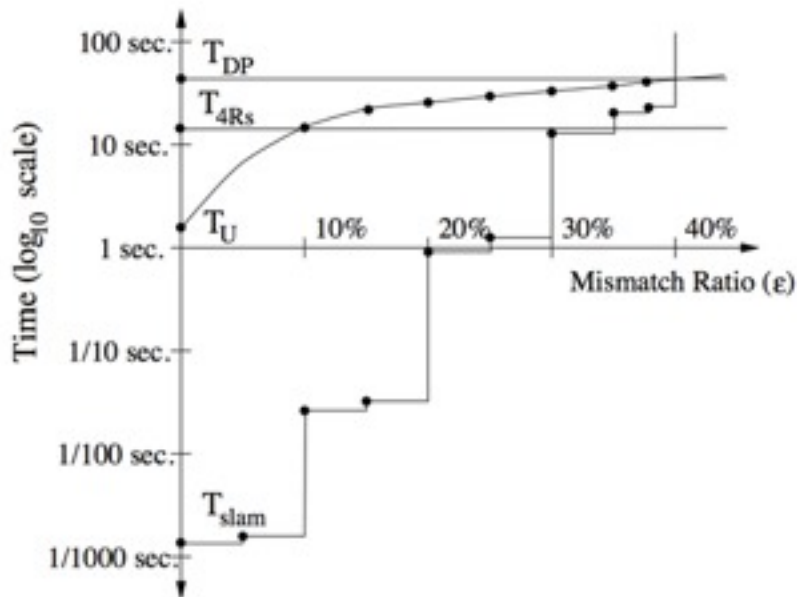


# The Result

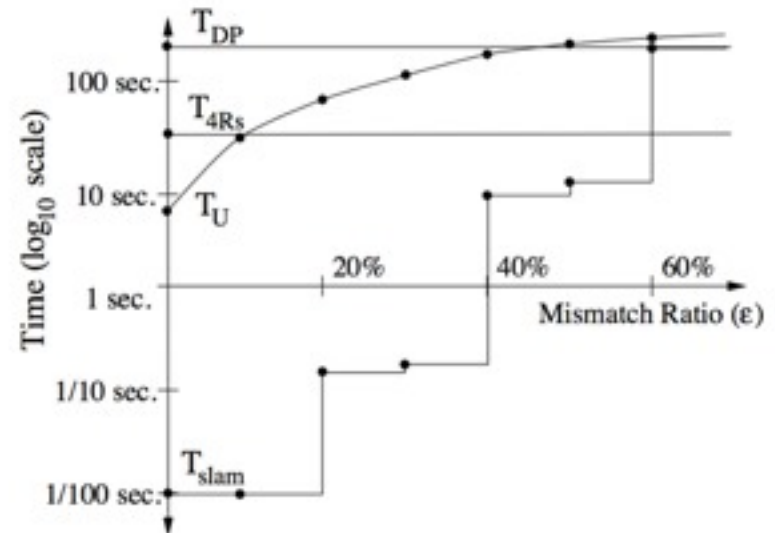
**Theorem:** Given

- (a)  $A$  is effectively Bernoulli,
- (b) a simple  $O(n)$  space, precomputed index of  $A$ , and
- (c) there are  $h$   $d$ -matches of a query  $Q$  to  $A$

then they can be found in  $O(d \cdot n^{\text{pow}(\epsilon)} \cdot \log n + pd \cdot h)$  expected-time.



$N = 1,000,000$  and  $|\Sigma| = 4$



$N = 4,000,000$  and  $|\Sigma| = 20$

**Figure 10:** Timing Plots for Queries of Length  $P = 80$

To my knowledge no one has improved  
on this in the last 20 years ! ?

*Algorithmica* 12, 4-5 (1994)

(submitted 1991 ! )



Les Treilles, May 1990