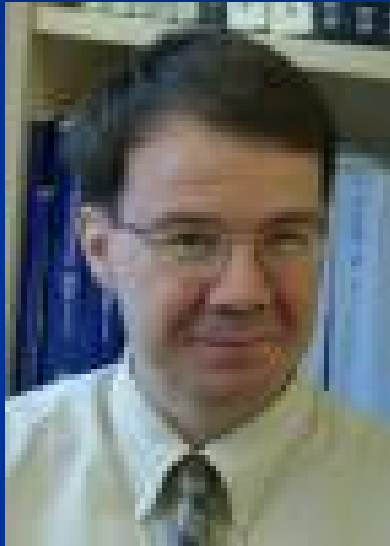


# Discovery of frequent episodes in event sequences using **WINEPI**

Maurizio Atzori  
University of Pisa

# The authors



Heikki Mannila



Hannu Toivonen



Inkeri Verkamo

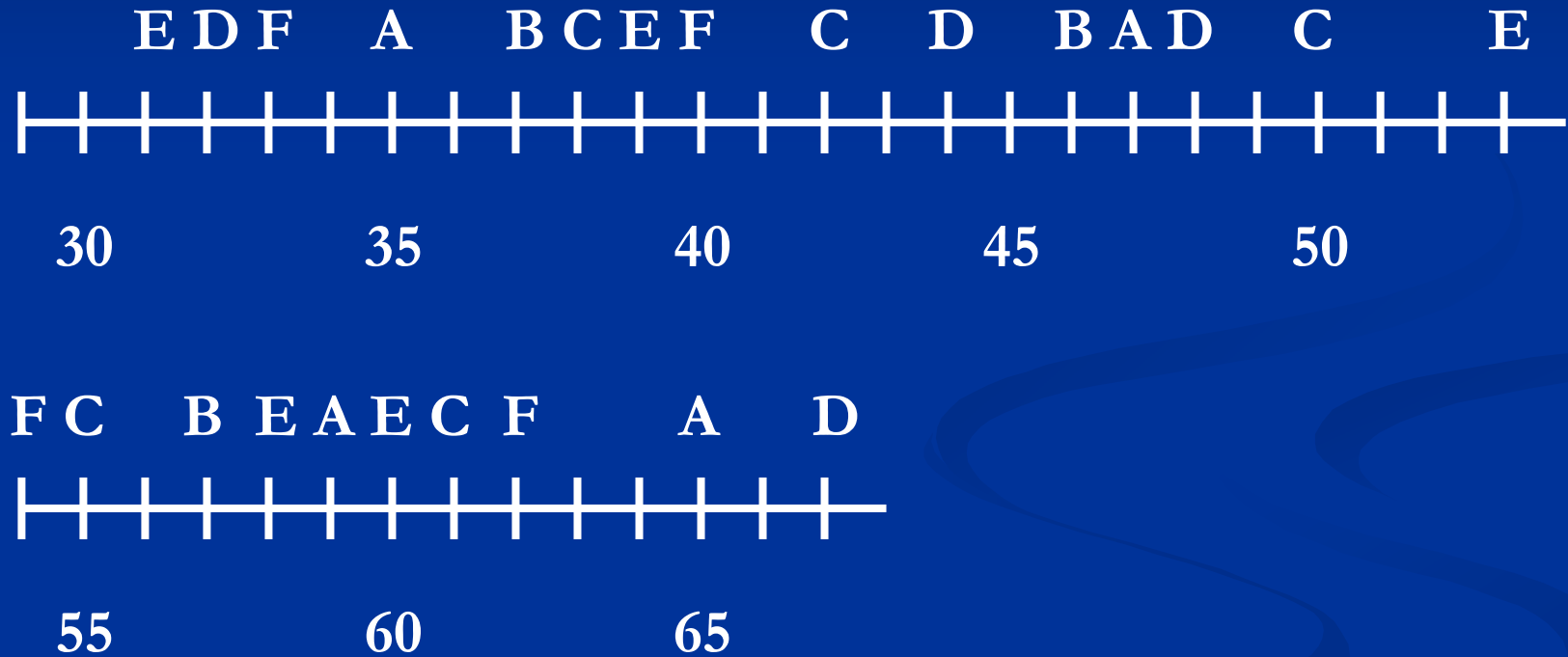
# Index

- Definitions and examples
  - Event sequences
  - Episodes
  - Windows
- Algorithms (WINEPI)
  - Based on finding “incrementally” frequent episodes
    - Idea taken from data mining association rules (apriori)
- Conclusions

# Definitions: event sequences

- An *event* is a pair  $(A, t)$  where
  - $A \in E$ , with  $E$  set of event types
  - $t$  is the *occurrence time* of the event
- An *event sequence  $s$  on  $E$*  is a triple  $(s, T_s, T_e)$  where:
  - $s = \langle (A_1, t_1), (A_2, t_2), \dots, (A_n, t_n) \rangle$
  - $t_i < t_{i+1} \quad \forall i = 1, \dots, n-1$

# Example of event sequence



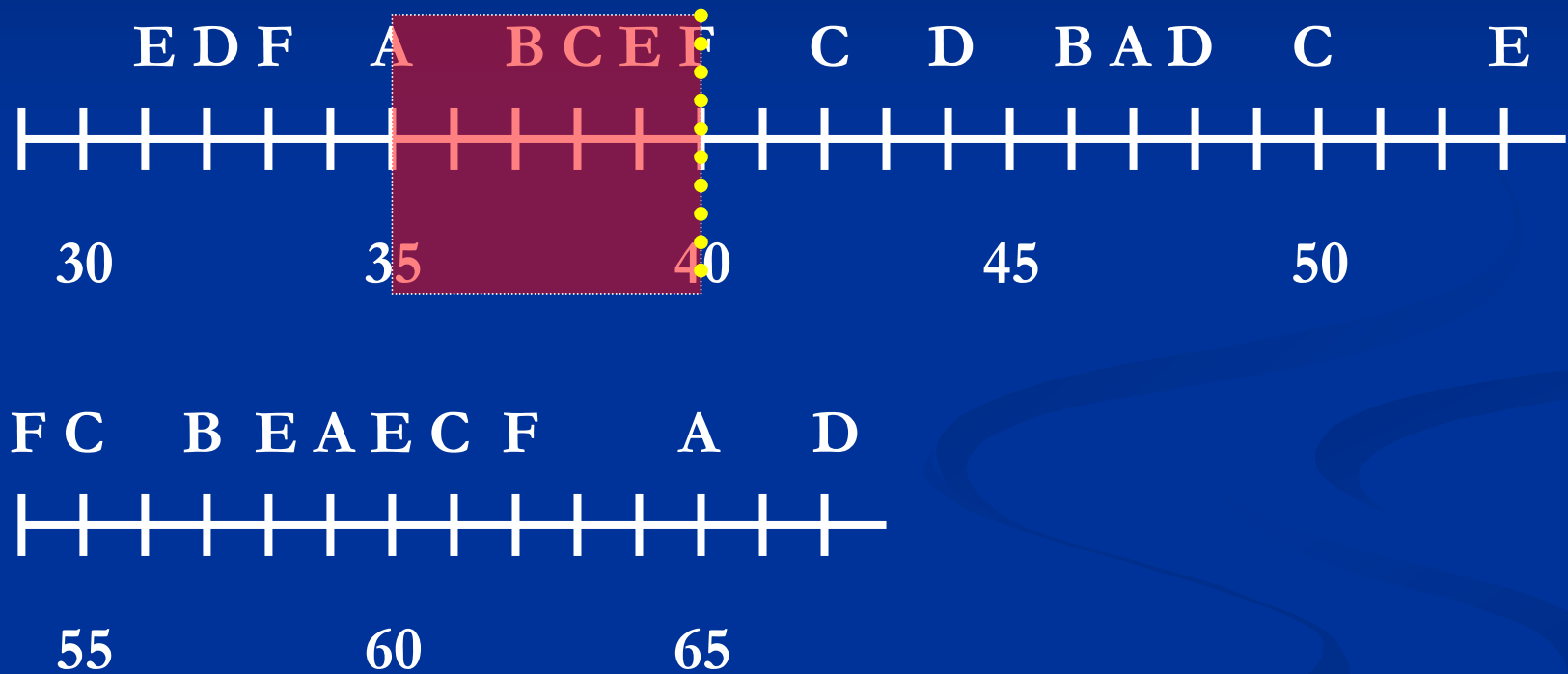
$s = \langle (E,31), (D,32), (F,33), (A,35), (B,37) \dots, (D,67) \rangle$

$s = (s, 29, 68)$

# Definitions: windows<sup>®</sup>

- A *window* on event sequence  $s=(s, T_s, T_e)$  is an event sequence  $w = (w, t_s, t_e)$  where:
  - $t_s < T_e, t_e > T_s$
  - $w$  consists of those pairs  $(A,t)$  from  $s$  where
    - $t_s \leq t < t_e$
- $t_e - t_s$  is called the *width* of the window  $w$ 
  - $\text{width}(w)$

# Example of window



$$w = \langle (A, 35), (B, 37), (C, 38), (E, 39) \rangle$$

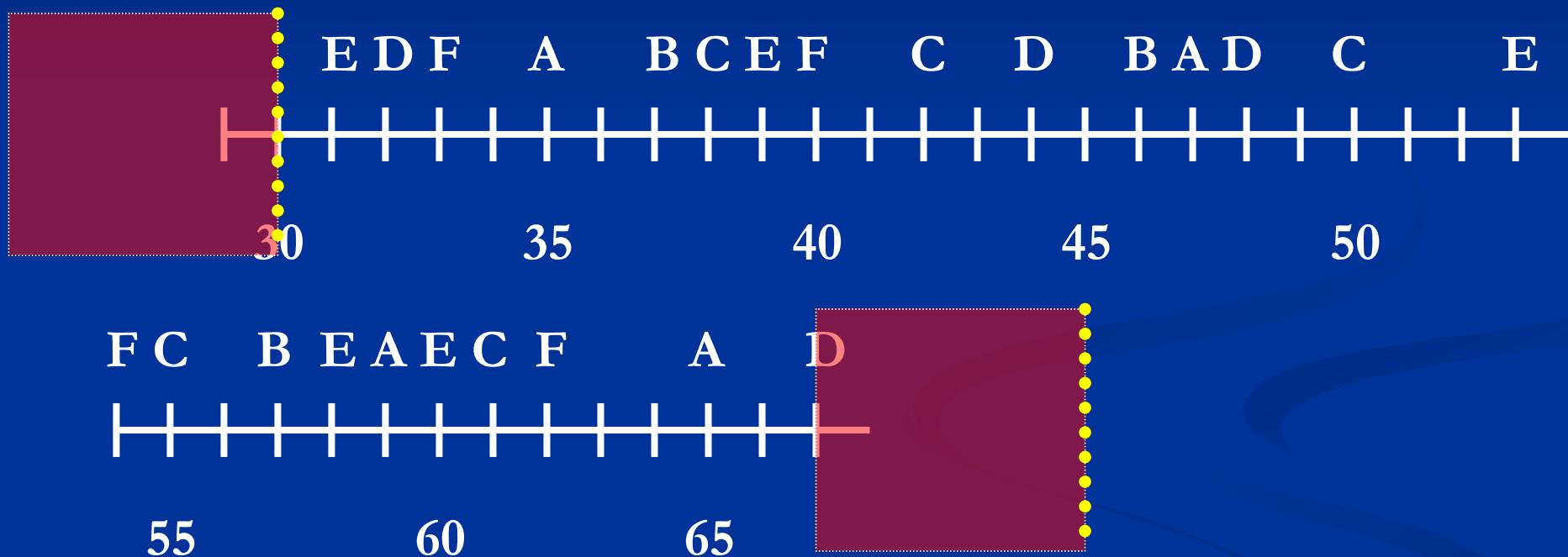
$$w = (w, 35, 40)$$

# Definitions: $W(.,.)$

- $W(s, win)$  is the set of all windows  $w$  on  $s$  such that  $width(w) = win$
- $\#W(s, win) = T_e - T_s + win - 1$
- In the previous figure,  $W(s, 5)$ :
  - The first window is  $(\phi, 25, 30)$
  - The last is  $(\langle(D, 67)\rangle, 67, 72)$



# Example of $W(.,.)$



■  $W(s, 5):$

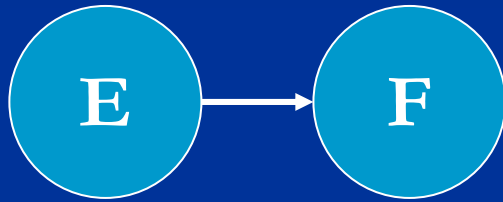
The first window is  $(\phi, 25, 30)$

The last is  $(\langle (D, 67) \rangle, 67, 72)$

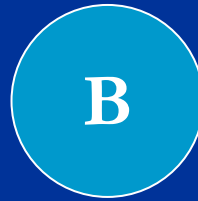
# Definitions: episodes

- Informally:
  - An *episode* is a partially ordered collection of events occurring together
- Formally:
  - An episode  $\alpha$  is a triple  $(V, \leq, g)$  where:
    - $V$  is a set of nodes,  $\leq$  is a partial order on  $V$
    - $g: V \rightarrow E$  is a mapping associating each node with an event type
  - The events in  $g(V)$  have to occur in the order described by  $\leq$
  - The size of  $\alpha$ , denoted by  $|\alpha|$ , is  $|V|$

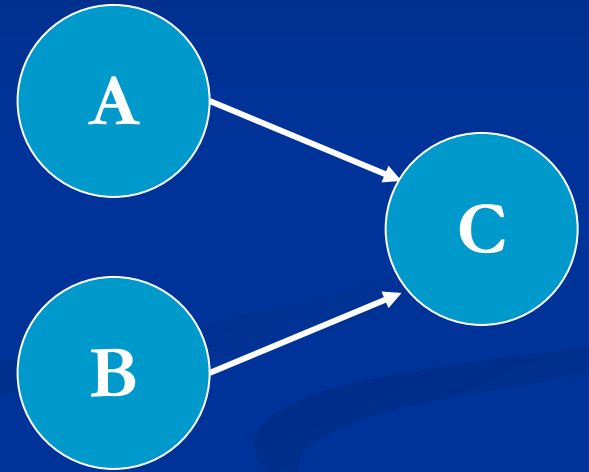
# Example of episodes



$\alpha$



$\beta$



$\gamma$

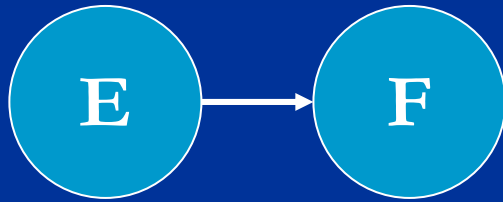
# Some kinds of episodes

- Parallel episodes
  - The partial order  $\leq$  is trivial
- Serial episodes
  - The  $\leq$  relation is a total order
- Injective episodes
  - The mapping  $g$  is an injection (i.e. no event type occurs twice in the episode)

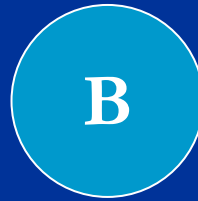
# Definitions: subepisodes

- An episode  $\beta = (V', \leq', g')$  is a *subepisode* of  $\alpha = (V, \leq, g)$ , denoted by  $\beta \leq \alpha$ , if
  - Exist an injective mapping  $f: V' \rightarrow V$  such that
    - $g'(v) = g(f(v)) \quad \forall v \in V'$
    - $\forall v, w \in V'. \quad v \leq' w \Rightarrow f(v) \leq f(w)$
- Informally:
  - $\beta \leq \alpha \Leftrightarrow$  the graph representing  $\beta$  is a subgraph of  $\alpha$

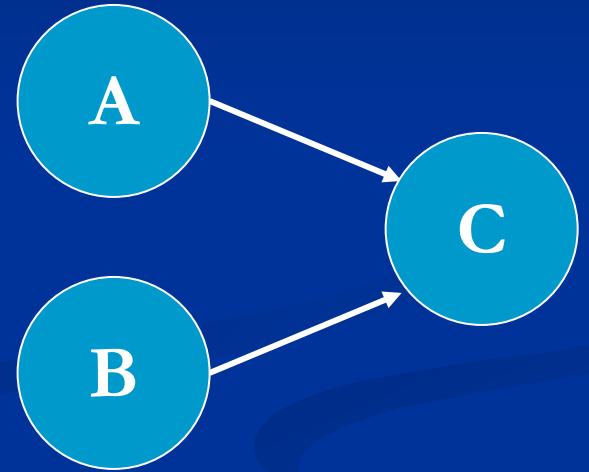
# Example of episodes



$\alpha$



$\beta$

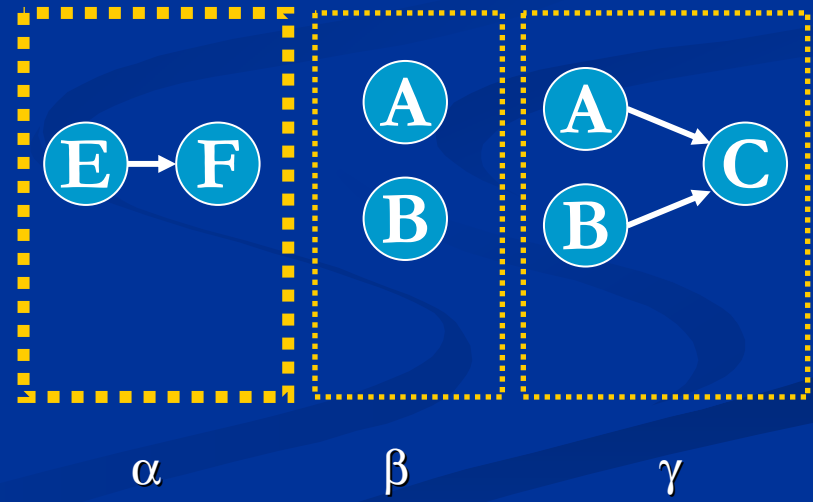
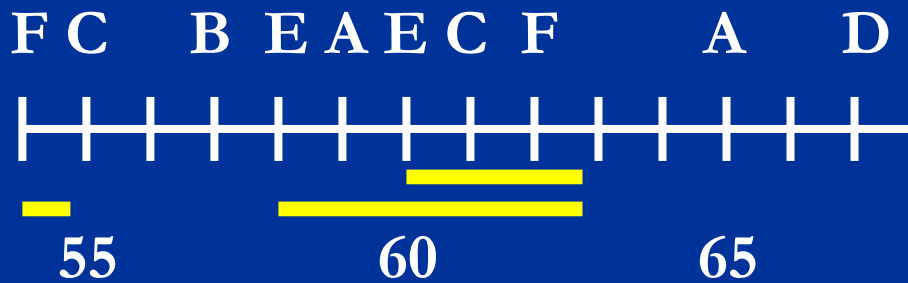
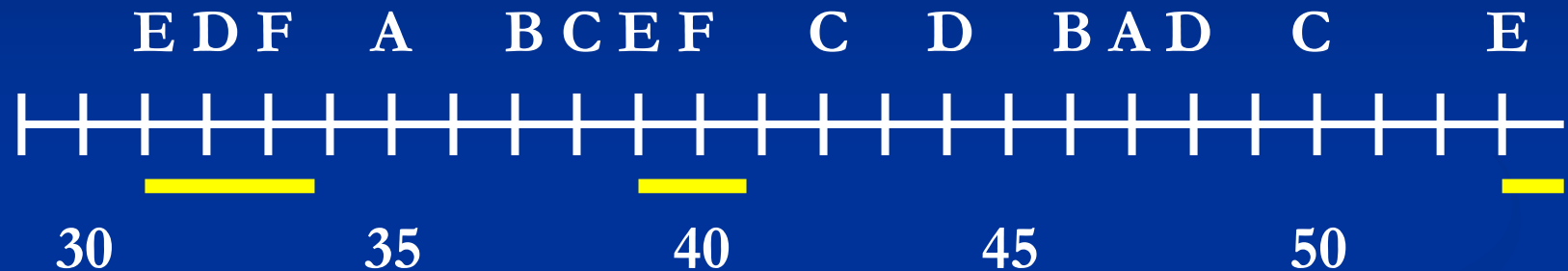


$\gamma$

# Definitions: occurring

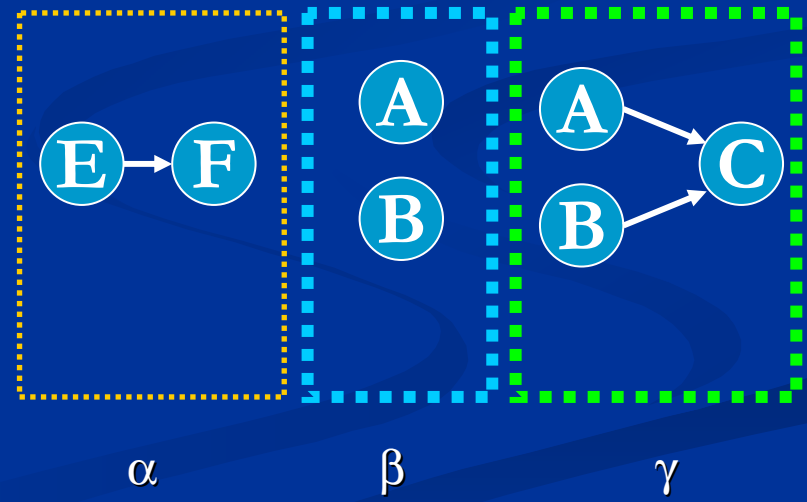
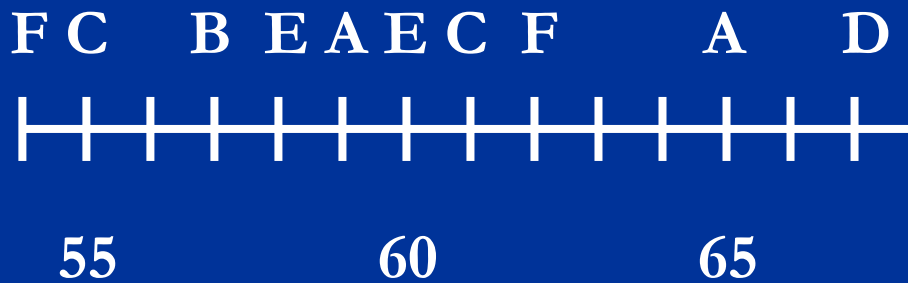
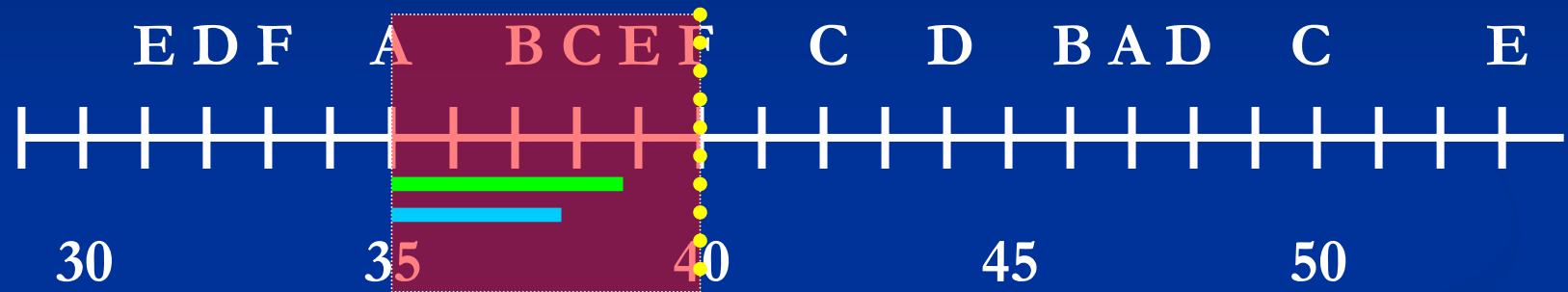
- An episode  $\alpha = (V, \leq, g)$  *occurs* in an event sequence  $s = \langle (A_1, t_1), (A_2, t_2), \dots, (A_n, t_n) \rangle$  if:
  - $\exists$  injective mapping  $h: V \rightarrow \{1, \dots, n\}$ 
    - $g(x) = A_{h(x)} \quad \forall x \in V$
    - $x \neq y$  and  $x \leq y \Rightarrow t_{h(x)} < t_{h(y)}$

# Example of occurring episodes in an event sequence





# Example of occurring episodes in a window



# Definitions: frequency

- The *frequency* of an episode is the fraction of windows in which the episode occurs

$$\text{fr}(\alpha, s, \text{win}) = \frac{|\{w \in W(s, \text{win}) \mid \alpha \text{ occurs in } w\}|}{|W(s, \text{win})|}$$

Given a *frequency threshold*  $\text{min\_fr}$ , an episode  $\alpha$  is *frequent* if  $\text{fr}(\alpha, s, \text{win}) \geq \text{min\_fr}$

# The task

- We are interested to discover all frequent episodes from a given class  $\mathcal{E}$  of episodes
  - The class could be, e.g., all parallel episodes or all serial episodes
- $F(s, \text{win}, \text{min\_fr})$  is the collection of frequent episodes with respect to  $s$ ,  $\text{win}$  and  $\text{min\_fr}$

# Algorithm 1: finding rules

**Input:** A set  $E$  of event types, an event sequence  $s$  over  $E$ , a set  $\mathcal{e}$  of episodes, a window width  $\text{win}$ , a frequency threshold  $\text{min-fr}$ , and a confidence threshold  $\text{min-conf}$ .

**Output:** The episode rules that hold in  $s$  with respect to  $\text{win}$ ,  $\text{min-fr}$ , and  $\text{min-conf}$ .

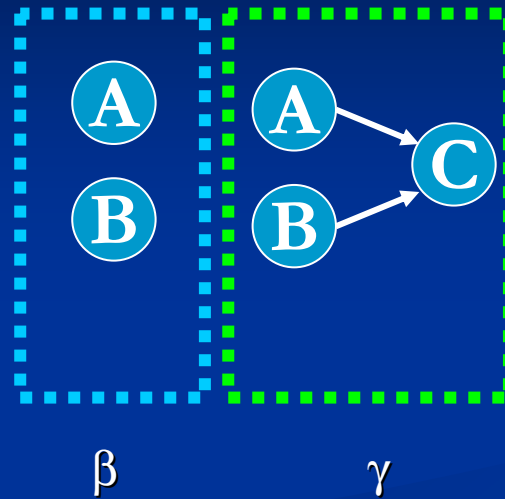
**Method:**

```
1. /* Find frequent episodes (Algorithm 2): */
2. compute  $F(s; \text{win}; \text{min-fr})$ ;
3. /* Generate rules: */
4. for all  $\alpha \in F(s; \text{win}; \text{min-fr})$  do
5.     for all  $\beta < \alpha$  do
6.         if  $\text{fr}(\alpha)/\text{fr}(\beta) \geq \text{min-conf}$  then
7.             output the rule  $\beta \rightarrow \alpha$  and the confidence
                $\text{fr}(\alpha)/\text{fr}(\beta)$ ;
```

# Example of rule

- Given:

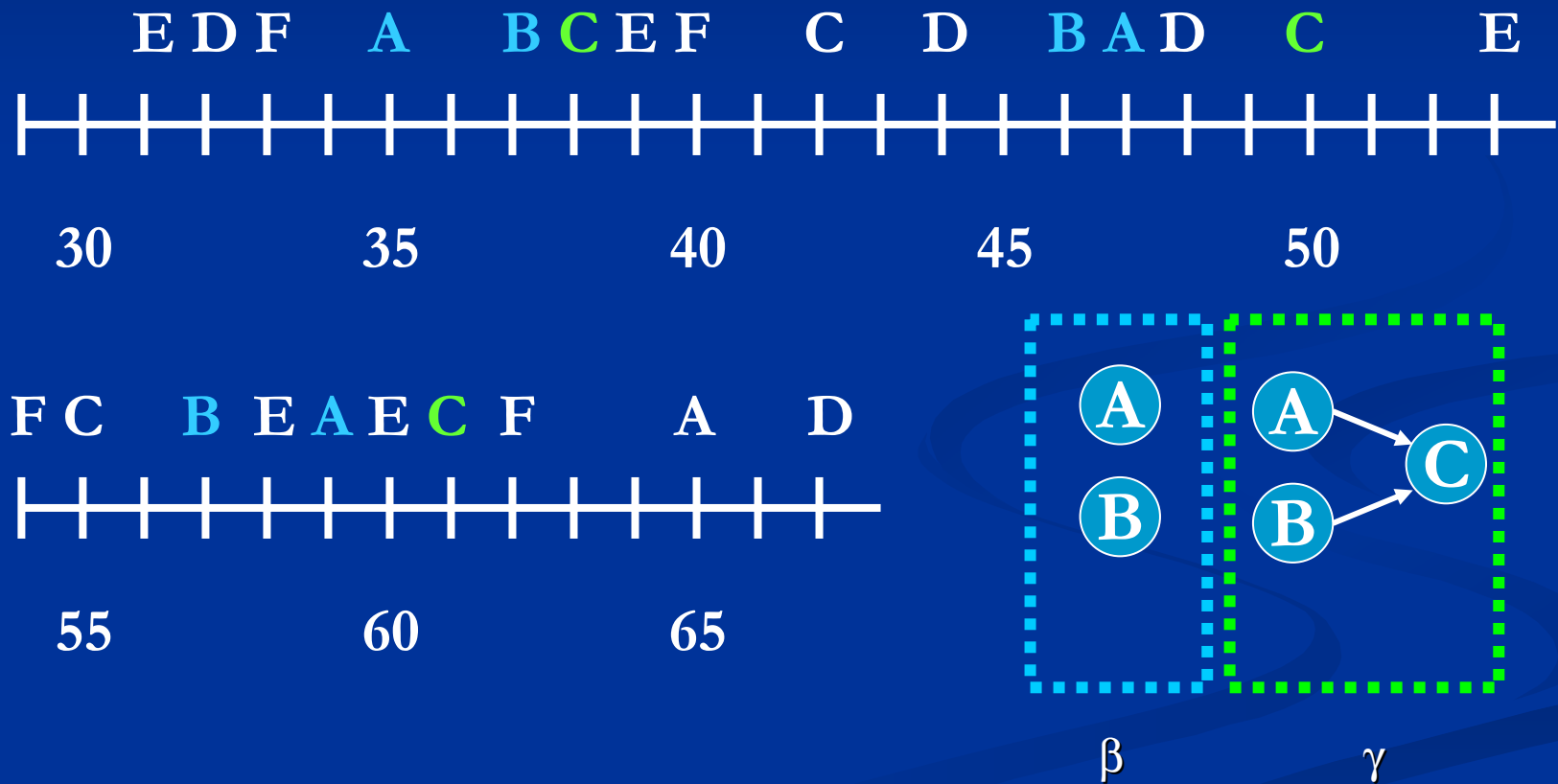
- $\beta < \gamma$
- $\text{fr}(\beta) = 4.2 \%$
- $\text{fr}(\gamma) = 4.0 \%$



- Then:

- If we have A and B in any order then there is a chance of 95% ( $4.0/4.2$ ) that C follows in the same window
- $A \text{ and } B \rightarrow C \text{ (95\%)}$

# Example of rule (2)



# The idea of WINEPI

- If an episode  $\alpha$  is frequent in an event sequence  $s$ , then all subepisodes  $\beta \leq \alpha$  are frequent
- This is used during the generation of candidate episodes, from small to big episodes

# Algorithm 2: compute $F(.,.)$

**Input:** A set  $E$  of event types, an event sequence  $s$  over  $E$ , set  $\varepsilon$  of episodes, a window width  $\text{win}$ , and a frequency threshold  $\text{min-fr}$ .

**Output:** The collection  $F(s; \text{win}; \text{min-fr})$  of frequent episodes.

**Method:**

1. compute  $C_1 := \{\alpha \in \varepsilon \mid |\alpha| = 1\};$
2.  $l := 1;$
3. while  $C_l \neq \emptyset$  do
4.     */\* Database pass (Algorithms 4 and 5): \*/*
5.     compute  $F_l := \{\alpha \in C_l \mid \text{fr}(\alpha; s; \text{win}) \geq \text{min-fr}\};$
6.      $l := l + 1;$
7.     */\* Candidate generation (Algorithm 3): \*/*
8.     compute  $C_l := \{\alpha \in \varepsilon \mid |\alpha| = l \text{ and for all } \beta \in \varepsilon \text{ such that } \beta < \alpha \text{ and } |\beta| < l \text{ we have } \beta \in F_{|\beta|}\};$
9. for all  $l$  do output  $F_l$  ;



# Generation of candidate episodes

- $C_l$  is the set containing candidate episodes of length  $l$ 
  - Generation step
- $F_l$  is the set containing all frequent episodes of length  $l$ 
  - “Pruning” step

# Algorithm 3: candidate episodes generation

**Input:** A sorted array  $F_l$  of frequent parallel episodes of size  $l$ .

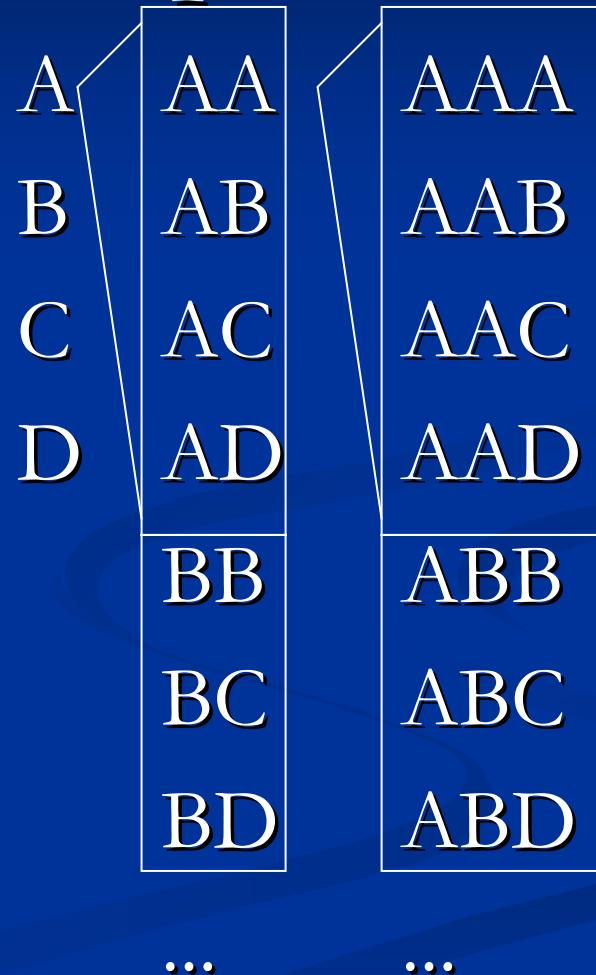
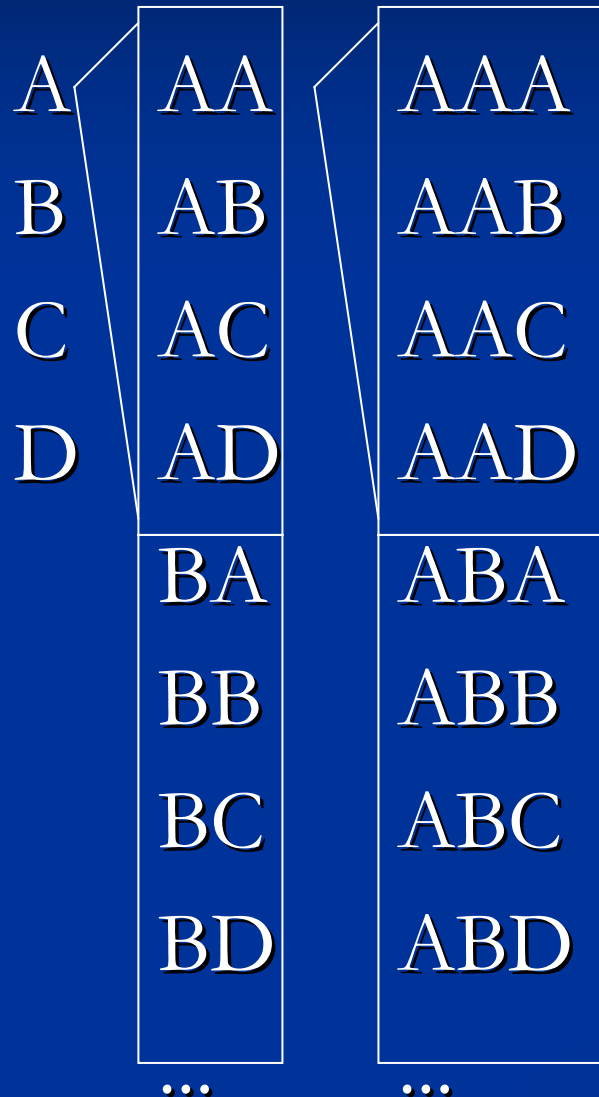
**Output:** A sorted array  $(C_{l+1})$  of candidate parallel episodes of size  $l + 1$ .

# Algorithm 3: two steps

- Build a potential candidate  $\alpha$  as a combination of 2 episodes from the same block
- Build and test subepisodes  $\beta$  that do not contain  $\alpha[y]$

# Blocks

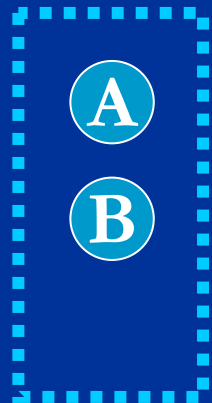
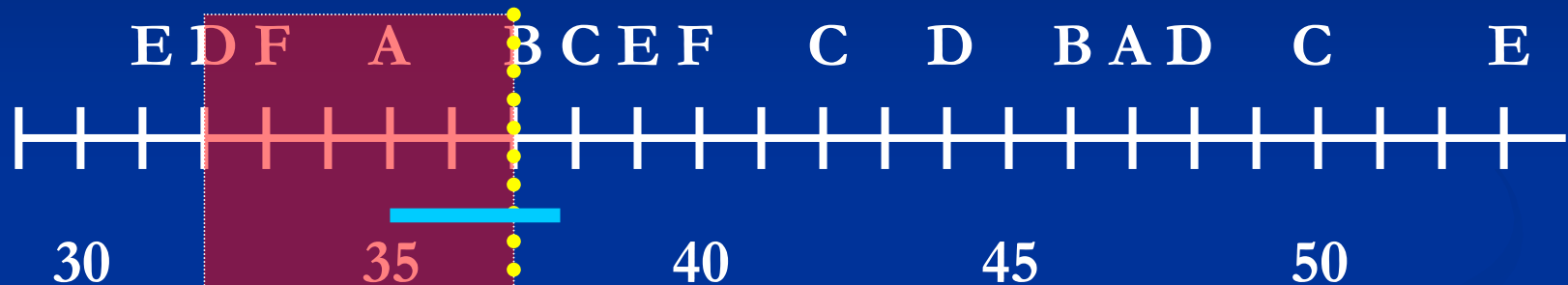
serials / parallel episodes



# Algorithm 4: recognizing parallel episodes

1. Initialization
2. Recognition
  1. Bring in new events to the window
  2. Drop out old events from the window

# Algorithm 4: example



D, F, A.count = 1

**Bring in**

B.count := B.count + 1

Look at Contains(B, 1)

**Drop out**

Look at Contains(D, 1)

D.count := D.count - 1

If  $\alpha.\text{event\_count} = |\alpha|$   
 $\alpha.\text{inwin} = \text{start}$

If  $\alpha.\text{event\_count} = |\alpha|$

$\alpha.\text{freq\_count} +=$   
 $\text{start} - \alpha.\text{inwin}$

# Algorithm 5: recognizing serial episodes

- The basic idea is to have an automaton for each episode

- Data structures:

$(\alpha, x)$  episode  $\alpha$  is waiting for  $x$ th event

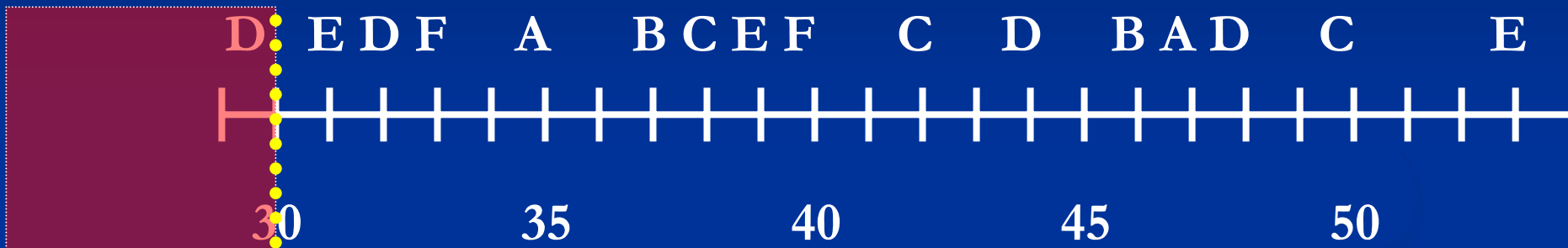
$\text{Wait}(A)$  set of episodes waiting for  $A$

$\alpha.\text{initialized}[i]$  the latest initialization time of an automaton that has reached its  $i$ th state

$\text{transitions}$  set of transitions  $(\alpha, x, t)$

episode  $\alpha$  got its  $x$ th event at time  $t$

# Algorithm 5: example



Look at  $\text{Wait}(D)$

transitions  $U = \{(\alpha, 1, \text{start} + \text{win} - 1)\}$



# Complexity

- Algorithm 3  $O(l^2 |F_1| \log |F_1|)$
- Algorithm 4  $O((n+l^2) |C|)$   
In practice,  $l \ll n$
- Algorithm 5  $O(n |C| 1)$

# Example of episodes

