# EPISODESUPPORT: A GLOBAL CONSTRAINT FOR MINING FREQUENT PATTERNS IN A LONG SEQUENCE OF EVENTS

Q. Cappart[1], **J. AOGA**[1], P. Schaus[1]

[1]UCLouvain — Belgium

UCL Université catholique de Louvain
icteam

CPAIOR 2018, Delft, The Netherlands, 26–29/06/2018

---

## CONTEXT

▷ This talk is about finding **frequent episode patterns** in a **long (time-stamped) sequence**

  ○ Very efficient dedicated algorithms exists (Minepi, Winepi, Emma,…)

    ○ They are not **flexible** and suffer for **memory problem**

▷ Motivation for CP:

  ○ finding **frequent (constrained) sub-sequences** is a related problem to **frequent (constrained) episode patterns** in a long (time-stamped) sequence

    ▫ constraint example: satisfying a regular constraint

  ○ CP-based method is the state-of-the-art for finding **frequent (constrained) sub-sequences** in a sequence database *[Aoga et al., ECMLPKDD'16; CPAIOR'17]*

---

## SPM PROBLEM

**Sub sequence**      **Sequence**

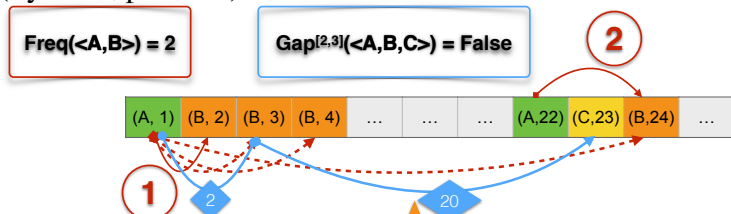| | | | | |
|---|---|---|---|---|
| **Client1** | Milk | **Coffee** | Sugar | Coffee | Sugar |
| **Client2** | Coffee | Milk | Coffee | Sugar | |
| **Client3** | Milk | Coffee | | | |
| **Client4** | Coffee | Sugar | Egg | | |

Sequence Database (SDB)

- Sequence : < Milk Coffee Sugar Coffee Sugar>
- Subsequence : <Coffee Sugar>
- Freq (<Coffee Sugar>) = 3

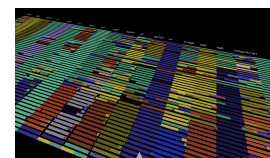**Problem : Find all subsequences with support ≥ Given Threshold**

---

## (CONSTRAINT-BASED) EPISODE MINING problem

(Symbol, position)

Freq(<A,B>) = 2    Gap$^{[2,3]}$(<A,B,C>) = False    **2**

| (A, 1) | (B, 2) | (B, 3) | (B, 4) | … | … | … | (A,22) | (C,23) | (B,24) | … |

**1**  2  20

**Problem :** Find **All Episodes** wrt. user-defined constraints (e.g. support ≥ Given Threshold)

☑ Applications



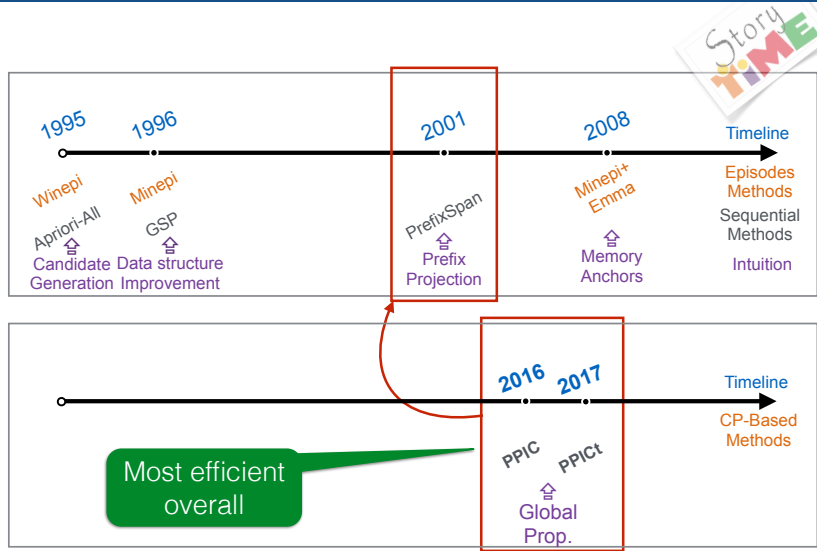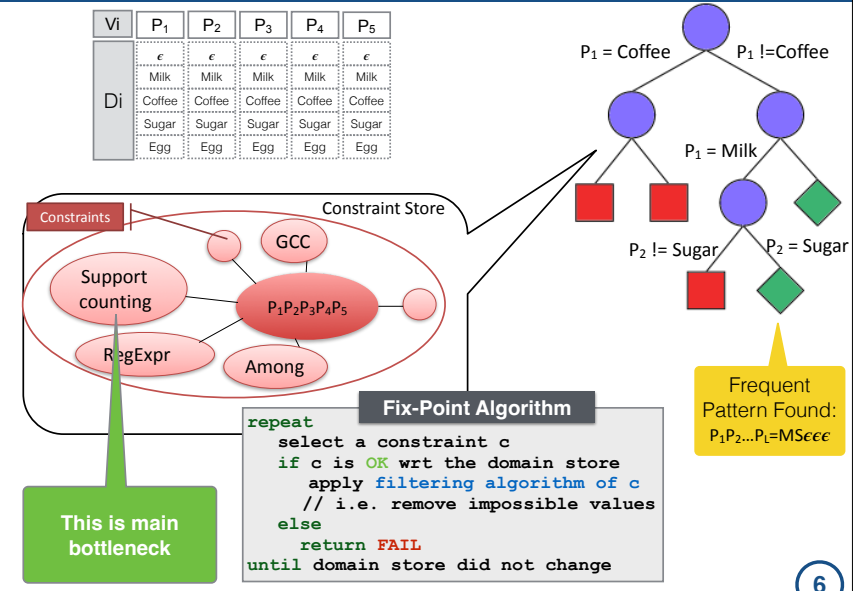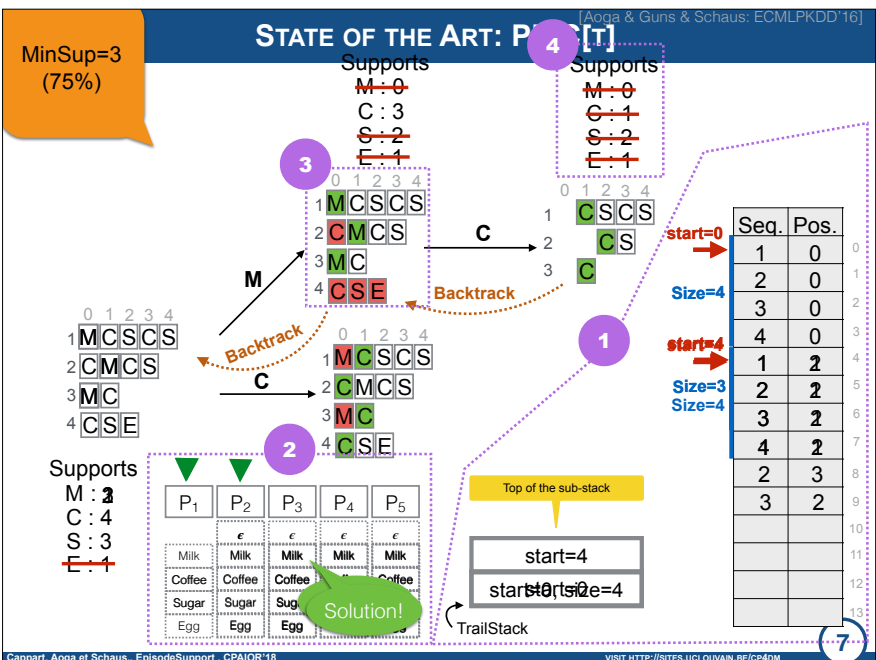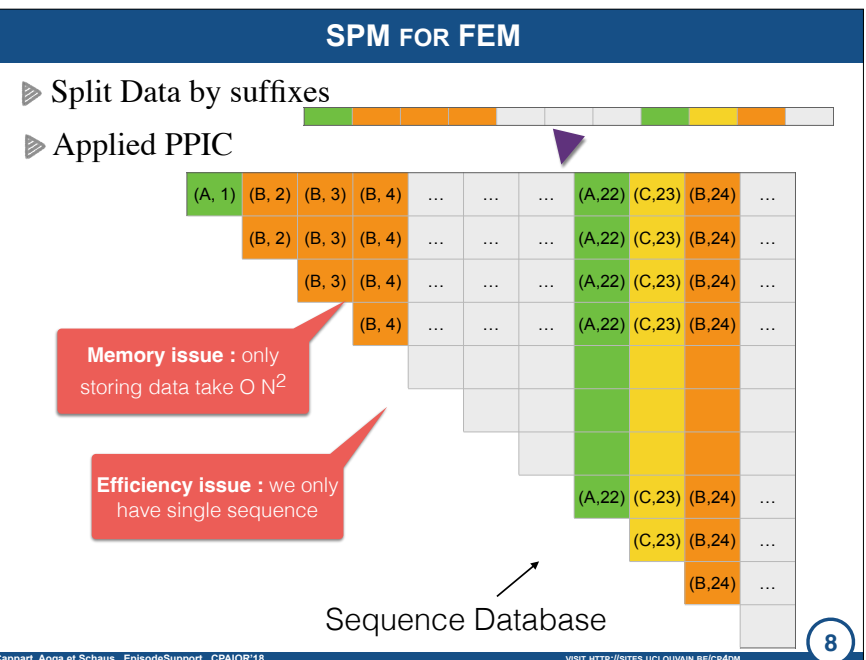DNA Sequence    Smartphone lifelogging    Stock market

# Related Work



1995  1996  2001  2008  Timeline

Winepi  Minepi  PrefixSpan  Minepi+ Emma  Episodes Methods

Apriori-All  GSP  Sequential Methods

Candidate Generation  Data structure Improvement  Prefix Projection  Memory Anchors  Intuition

2016  2017  Timeline

PPIC  PPICt  CP-Based Methods

Global Prop.

Most efficient overall

Story TIME

---

# CP : Filtering + DFSearch



| Vi | P₁ | P₂ | P₃ | P₄ | P₅ |
|---|---|---|---|---|---|
| Di | ε | ε | ε | ε | ε |
| | Milk | Milk | Milk | Milk | Milk |
| | Coffee | Coffee | Coffee | Coffee | Coffee |
| | Sugar | Sugar | Sugar | Sugar | Sugar |
| | Egg | Egg | Egg | Egg | Egg |

Constraints  Constraint Store

GCC

Support counting

$P_1 P_2 P_3 P_4 P_5$

RegExpr  Among

This is main bottleneck

$P_1$ = Coffee   $P_1$ !=Coffee

$P_1$ = Milk

$P_2$ != Sugar   $P_2$ = Sugar

Frequent Pattern Found: $P_1 P_2 ... P_L = MS\epsilon\epsilon\epsilon$

**Fix-Point Algorithm**

```
repeat
    select a constraint c
    if c is OK wrt the domain store
        apply filtering algorithm of c
        // i.e. remove impossible values
    else
        return FAIL
until domain store did not change
```

---

# STATE OF THE ART: PPIC [T]

[Aoga & Guns & Schaus: ECMLPKDD'16]

MinSup=3 (75%)



Supports
M : 0
C : 3
S : 2
E : 1

Supports
M : 0
C : 1
S : 2
E : 1

| | 0 1 2 3 4 |
|---|---|
| 1 | M C S C S |
| 2 | C M C S |
| 3 | M C |
| 4 | C S E |

C  start=0

Backtrack  M  Backtrack

| | 0 1 2 3 4 |
|---|---|
| 1 | M C S C S |
| 2 | C M C S |
| 3 | M C |
| 4 | C S E |

C

| | 0 1 2 3 4 |
|---|---|
| 1 | C S C S |
| 2 | C S |
| 3 | C |

start=4

| Seq. | Pos. |
|---|---|
| 1 | 0 |
| 2 | 0 |
| 3 | 0 |
| 4 | 0 |
| 1 | 2 |
| 2 | 2 |
| 3 | 2 |
| 4 | 2 |
| 2 | 3 |
| 3 | 2 |

Size=4
Size=3
Size=4

Supports
M : 3
C : 4
S : 3
E : 1

| P₁ | P₂ | P₃ | P₄ | P₅ |
|---|---|---|---|---|
| ε | ε | ε | ε | ε |
| Milk | Milk | Milk | Milk | Milk |
| Coffee | Coffee | Coffee | Coffee | Coffee |
| Sugar | Sugar | Sugar | Sugar | Sugar |
| Egg | Egg | Egg | Egg | Egg |

Top of the sub-stack

start=4
start=4  size=4

TrailStack

Solution!

---

# SPM FOR FEM

▷ Split Data by suffixes
▷ Applied PPIC



| (A, 1) | (B, 2) | (B, 3) | (B, 4) | ... | ... | ... | (A,22) | (C,23) | (B,24) | ... |
|---|---|---|---|---|---|---|---|---|---|---|
| | (B, 2) | (B, 3) | (B, 4) | ... | ... | ... | (A,22) | (C,23) | (B,24) | ... |
| | | (B, 3) | (B, 4) | ... | ... | ... | (A,22) | (C,23) | (B,24) | ... |
| | | | (B, 4) | ... | ... | ... | (A,22) | (C,23) | (B,24) | ... |
| | | | | | | | (A,22) | (C,23) | (B,24) | ... |
| | | | | | | | | (C,23) | (B,24) | ... |
| | | | | | | | | | (B,24) | ... |

**Memory issue :** only storing data take O $N^2$

**Efficiency issue :** we only have single sequence

Sequence Database

## Slide 9

**Goal:** Design new Approach for finding Episodes capturing the most common constraints (including syntax and time-related constraints)

☑ *Adapt* trailed-based data structure to efficiently overcome memory issue

☑ Take into account that we have a single sequence with algorithmic improvements

☑ Tackle time series data and time-related constraints

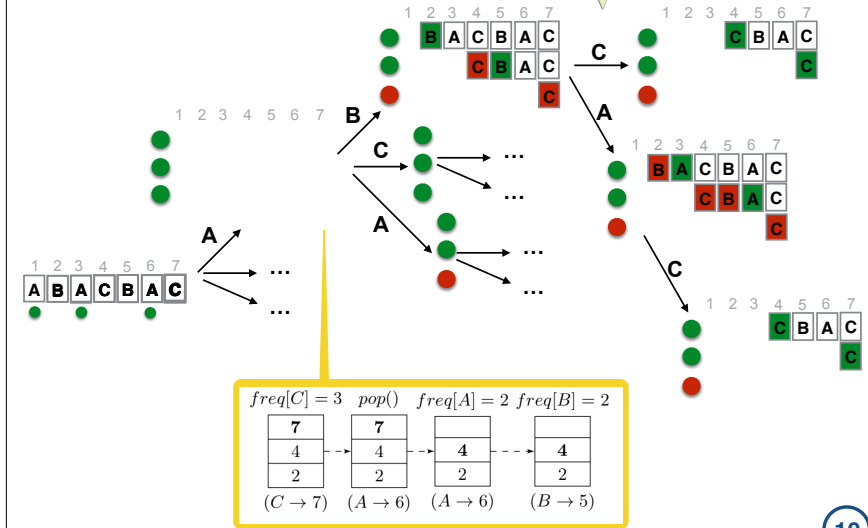☑ Show real application handling many other contraints: Regular/ Grammar, Gcc, Among, …

## Slide 10

## Slide 11 — EXPERIMENTS



**EXPERIMENTS**

OSCAR Scala
www.oscarlib.org

▷ EpisodeSuport (2versions — with/without time consideration)

https://bitbucket.org/projetsJOHN/episodesupport (also available in OscaR)

## Slide 12

### DFEM VS EPISODESUPPORT (MEMORY)

Time limit = 600s (10Minutes)

**Uniprot - Humain proteins dataset**
(2452 instances) - θ=5%; MaxSize=5



● DFEM    ◆ episodeSupport
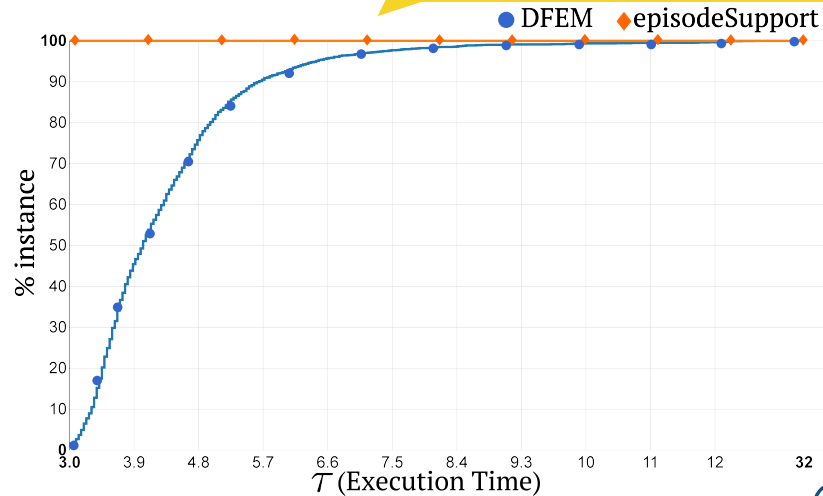
% instance vs $\mathcal{T}$ (Memory)

## Slide 13 — DFEM vs EpisodeSupport (Execution Time)

**Time limit = 600s (10Minutes)**

**Uniprot - Humain proteins dataset** (2452 instances) - θ=5%; MaxSize=5

Legend: ● DFEM  ◆ episodeSupport

Y-axis: % instance (0–100)
X-axis: $\mathcal{T}$ (Execution Time) — 3.0, 3.9, 4.8, 5.7, 6.6, 7.5, 8.4, 9.3, 10, 11, 12, 32
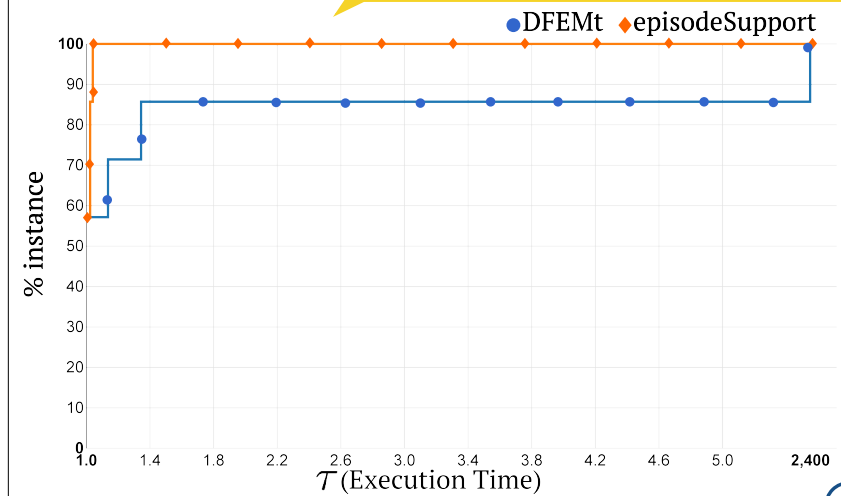
## Slide 14 — DFEMt vs EpisodeSupport (Execution Time)

**Time limit = 600s (10Minutes)**

**smartphone lifeloggins dataset** (21 instances) - θ=5%; MaxSize=5 gap[1s,1hour]

Legend: ● DFEMt  ◆ episodeSupport

Y-axis: % instance (0–100)
X-axis: $\mathcal{T}$ (Execution Time) — 1.0, 1.4, 1.8, 2.2, 2.6, 3.0, 3.4, 3.8, 4.2, 4.6, 5.0, 2,400

## Slide 15 — EpisodeSupport vs Minepi+ and Emma

Legend: MINEPI+  EMMA  EpisodeSupport

Our Method

Y-axis: Time/Average (log,s) — 0,00 / 2,25 / 4,50 / 6,75 / 9,00
X-axis: Datasets — USER3, USER2, USER1, USER5, USER4, USER8, USER6

## Slide 16 — Take-Away message

- Two versions of Global constraints (with/without time consideration)

- Efficiently split long sequence into small sequences for efficient memory usage

- Many kind of existing modules (in CP-Solvers) are reusable for free

- **Efficient memory using Trail-based backtracking aware data structure adaptation**

Thank you!