

SliceNStitch: Continuous CP Decomposition of Sparse Tensor Streams



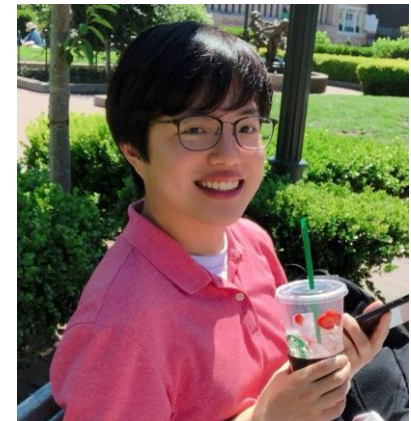
Taehyung Kwon*



Inkyu Park*



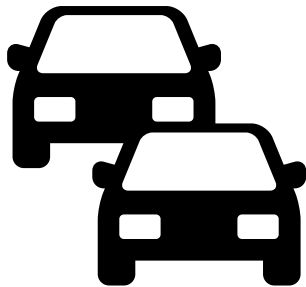
Dongjin Lee



Kijung Shin

Tensors are Everywhere

- Tensors are powerful tools for representing time-evolving multi-aspect data
- Multi-aspect data stream:
a sequence of timestamped M -tuples $\{(e_n = (i_1, \dots, i_{M-1}, v_n), t_n)\}_{n \in \mathbb{R}}$
 - i_1, \dots, i_{M-1} : non-time mode coordinates
 - v_n : value of the event
 - t_n : time when e_n occurs



(source, destination, 1)

Traffic history



(location, type, 1)

Crime history

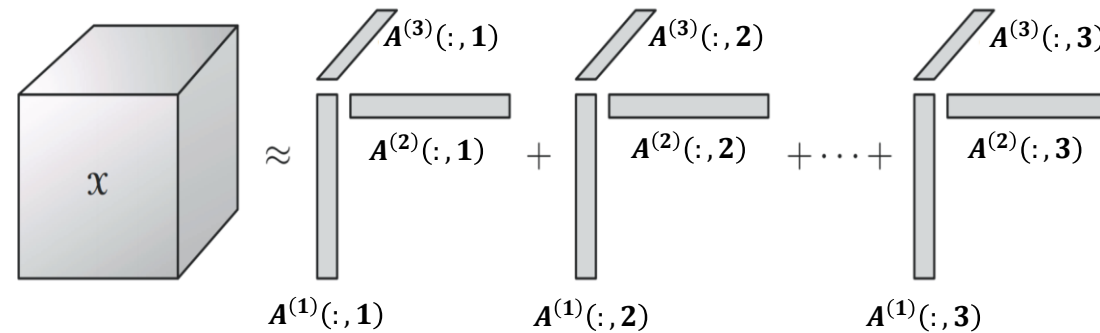


(user, product, color, quantity)

Purchase history

CANDECOMP/PARAFAC Decomposition (CPD)

- CPD gives a **low-rank approximation** $\tilde{\mathcal{X}}$ of the tensor \mathcal{X}
- Given an M -mode tensor $\mathcal{X} \in \mathbb{R}^{N_1 \times \dots \times N_M}$ and rank $R \in \mathbb{N}$,
 - $\mathcal{X} \approx \tilde{\mathcal{X}} \equiv \sum_{r=1}^R A^{(1)}(:, r) \circ \dots \circ A^{(M)}(:, r)$
 - where $A^{(1)}, \dots, A^{(M)}$: Factor matrices

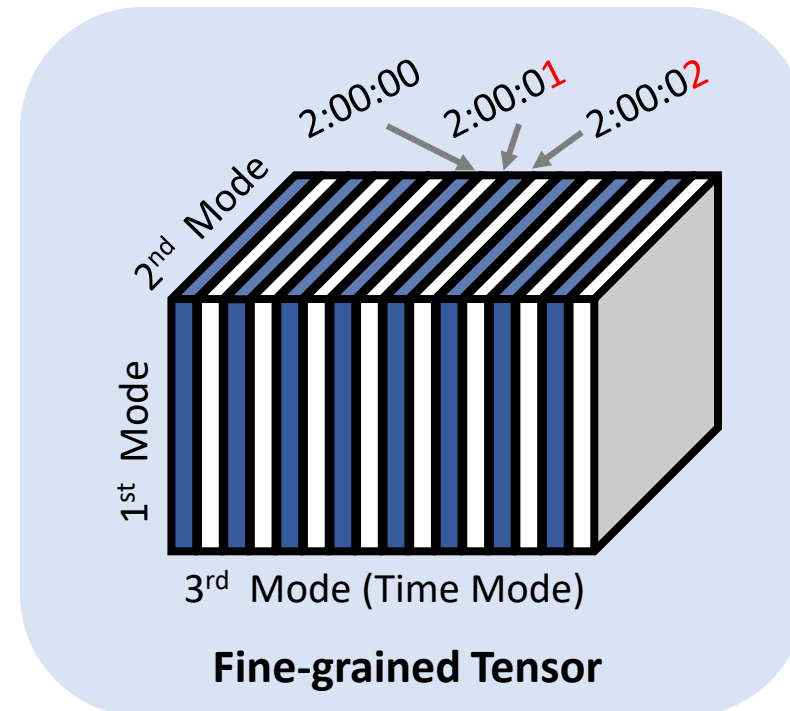
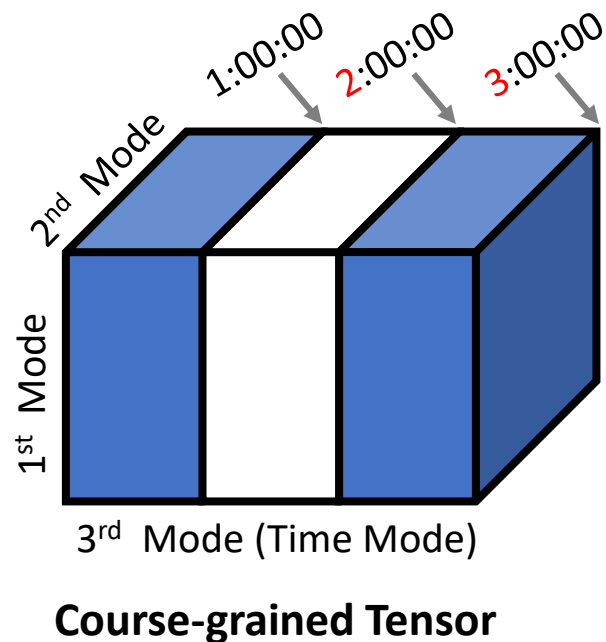


- Goal of CPD: factor matrices that minimize the difference between \mathcal{X} and $\tilde{\mathcal{X}}$

$$\min_{A^{(1)}, \dots, A^{(M)}} \|\mathcal{X} - \tilde{\mathcal{X}}\|_F$$

Limitation of Common Tensor Modeling

- Dynamic tensors **grow once per period**
 - Outputs of CPD are also updated once per period
- To perform CPD continuously for real-time application,
 - One can make the **granularity** of the time mode **extremely fine**



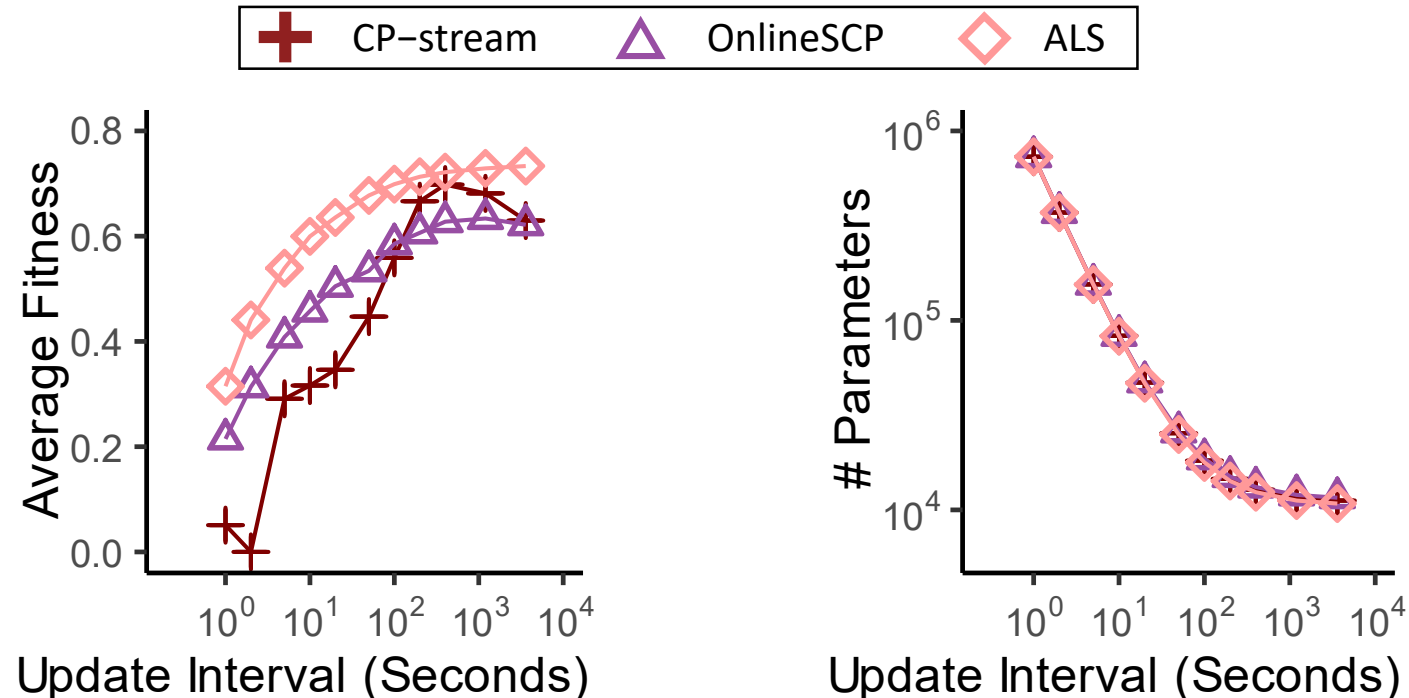
Limitation of Common Tensor Modeling

- Problems of fine-grained tensor modelings

- Degradation of fitness
- Increase the number of parameters

∴ Not suitable for real-time application!

| | Coarse-grained | Fine-grained |
|-----------------|----------------|--------------|
| Update Interval | Long (👎) | Short (👍) |
| Parameters | Few (👍) | Many (👎) |
| Fitness | High (👍) | Low (👎) |



Our Problem: Continuous CP Decomposition

- How can we continuously analyze multi-aspect data streams using CPD?
 - Given a multi-aspect data stream
 - Update its CP decomposition instantly in response to each new tuple in the stream
 - Without having to wait for the current period to end

Road Map

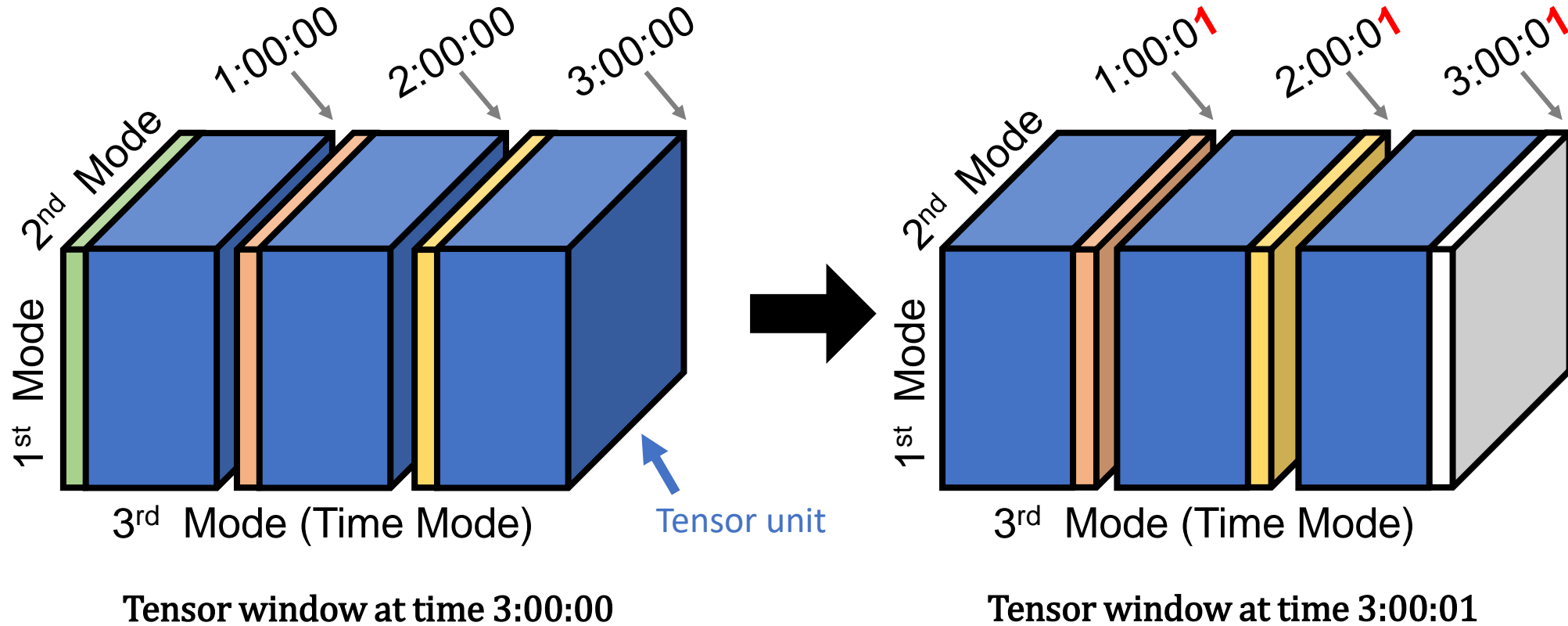
- Introduction
- Problem Definition
- **Data Model: Continuous Tensor Model** ←
- Optimization Algorithms: SliceNStitch
- Experiment Results
- Conclusions



Proposed Tensor Model

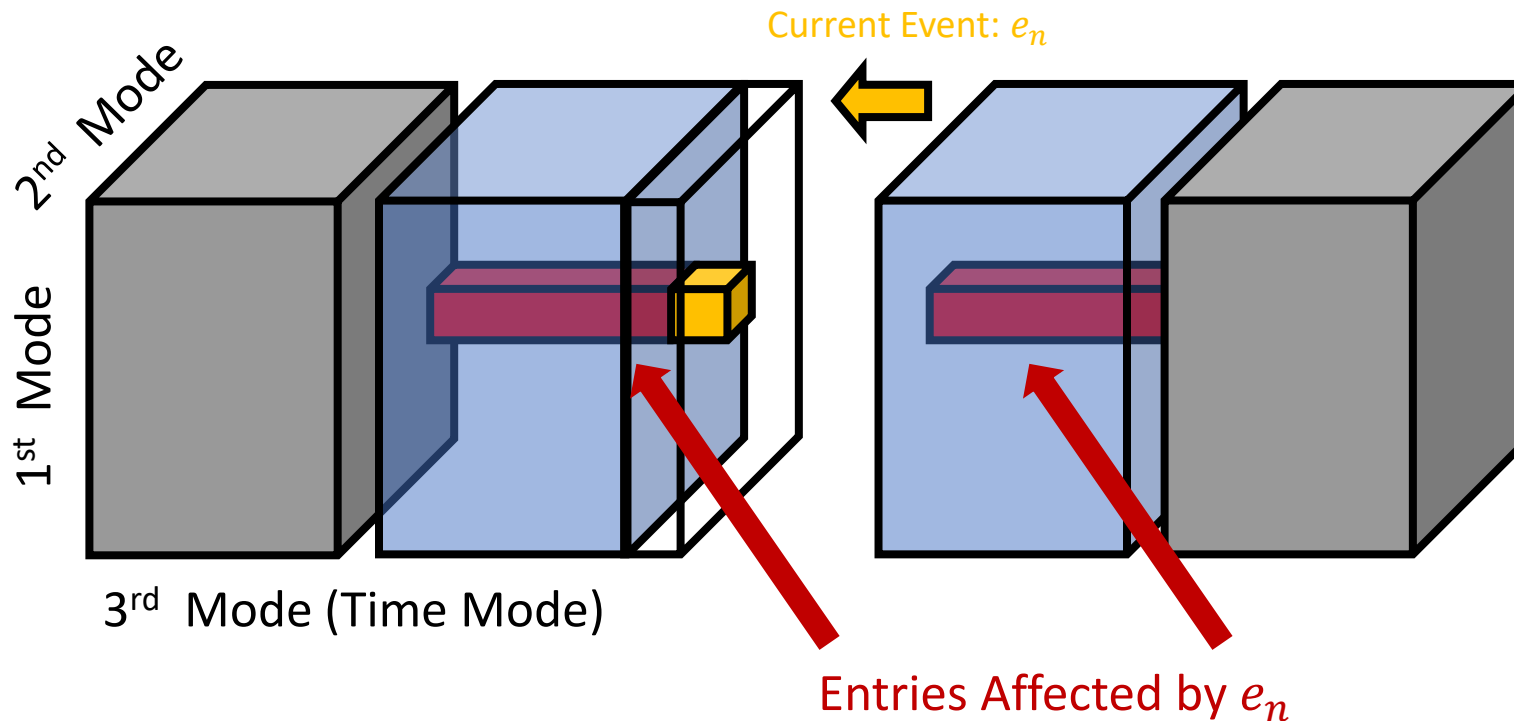
- Continuous Tensor Model

- The modeled tensor window and units evolve at each time



Event-driven Implementation

- For the tensor window \mathcal{X} ,
- a tuple $(e_n = (i_1, \dots, i_{M-1}, v), t_n)$ causes an event:



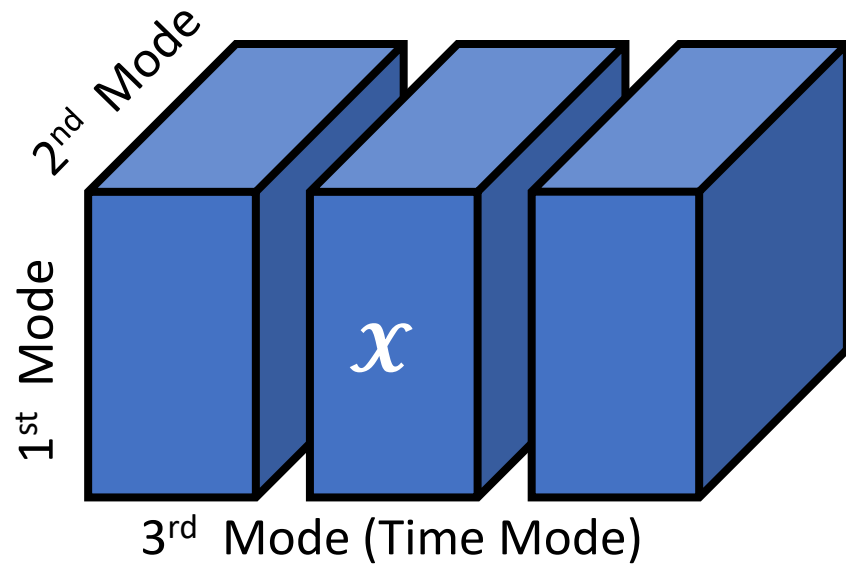
Road Map

- Introduction
- Problem Definition
- Data Model: Continuous Tensor Model
- **Optimization Algorithms: SliceNStitch** ←
- Experiment Results
- Conclusions

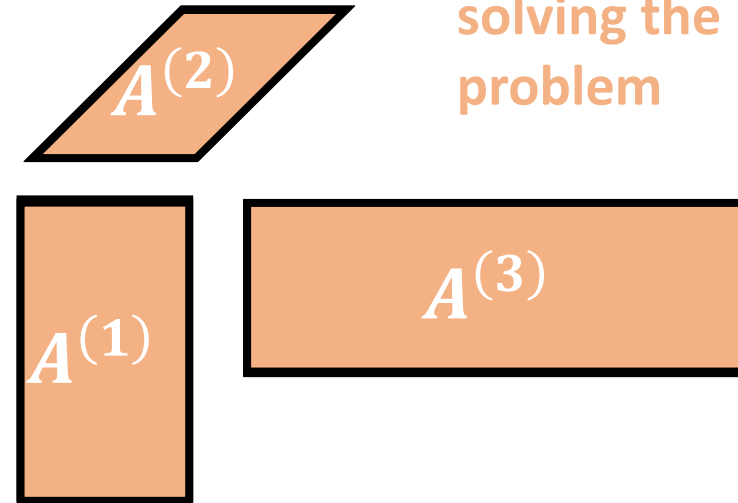


SliceNStitch-Matrix (SNS_{MAT})

Use all Non-zero Entries in \mathcal{X}



CP Decomposition



Update each matrix by solving the least square problem

Pros

High-quality solution

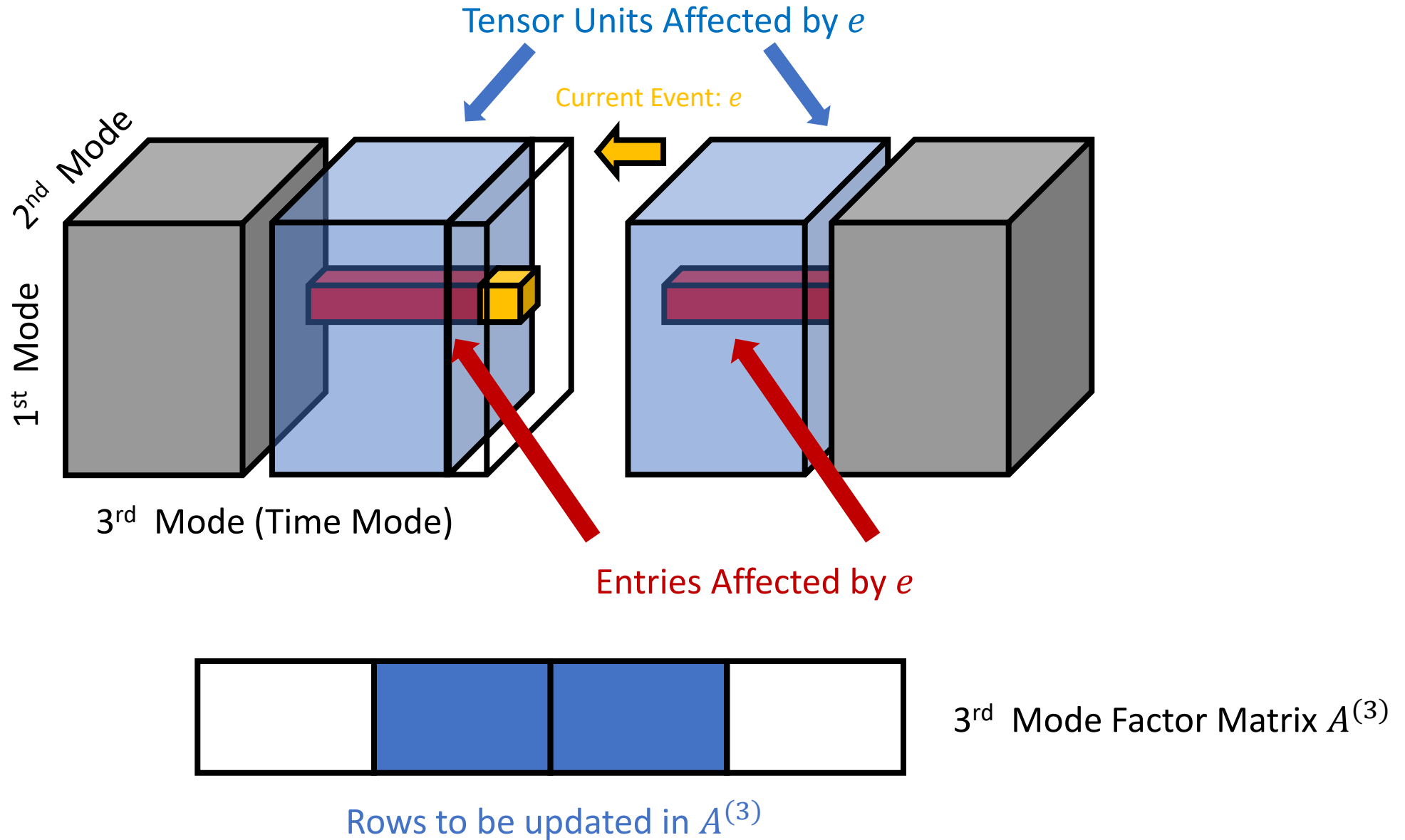


Cons

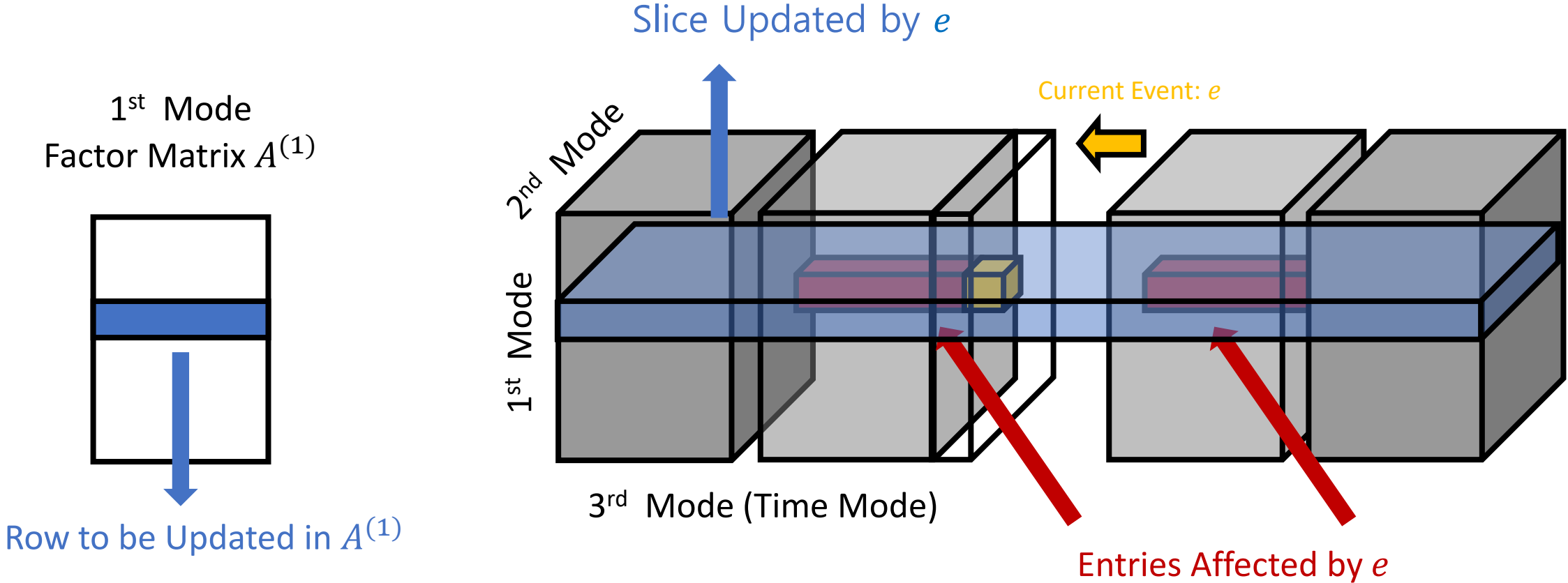
High computational cost



Common Outline of the Other Algorithms: Time Mode



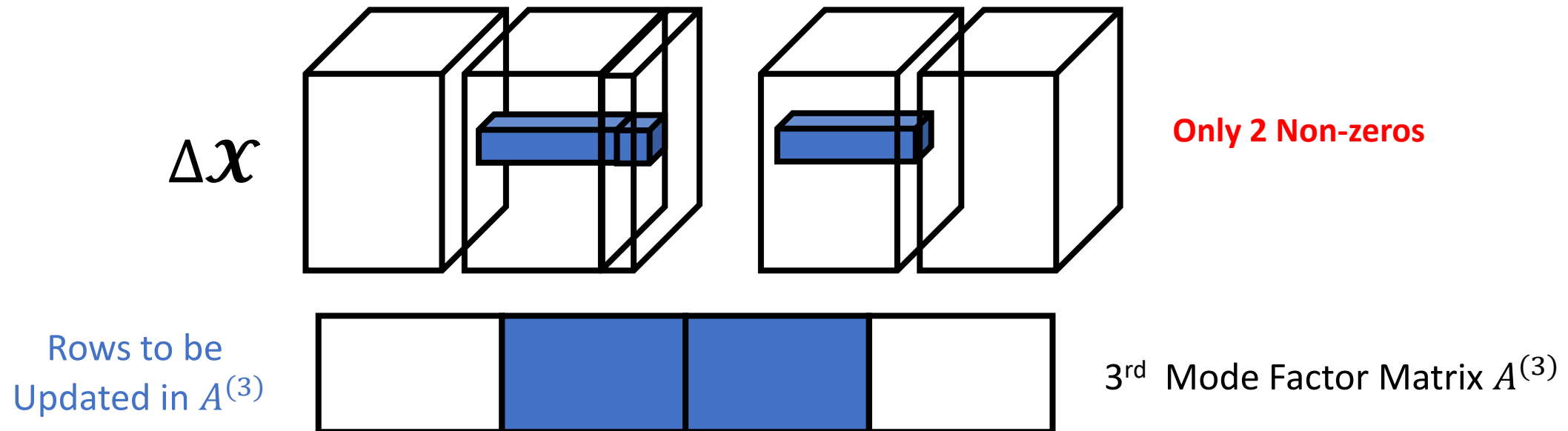
Common Outline of the Other Algorithms: Non-Time Mode



SliceNStitch-Vector: SNS_{VEC}

For the time mode,

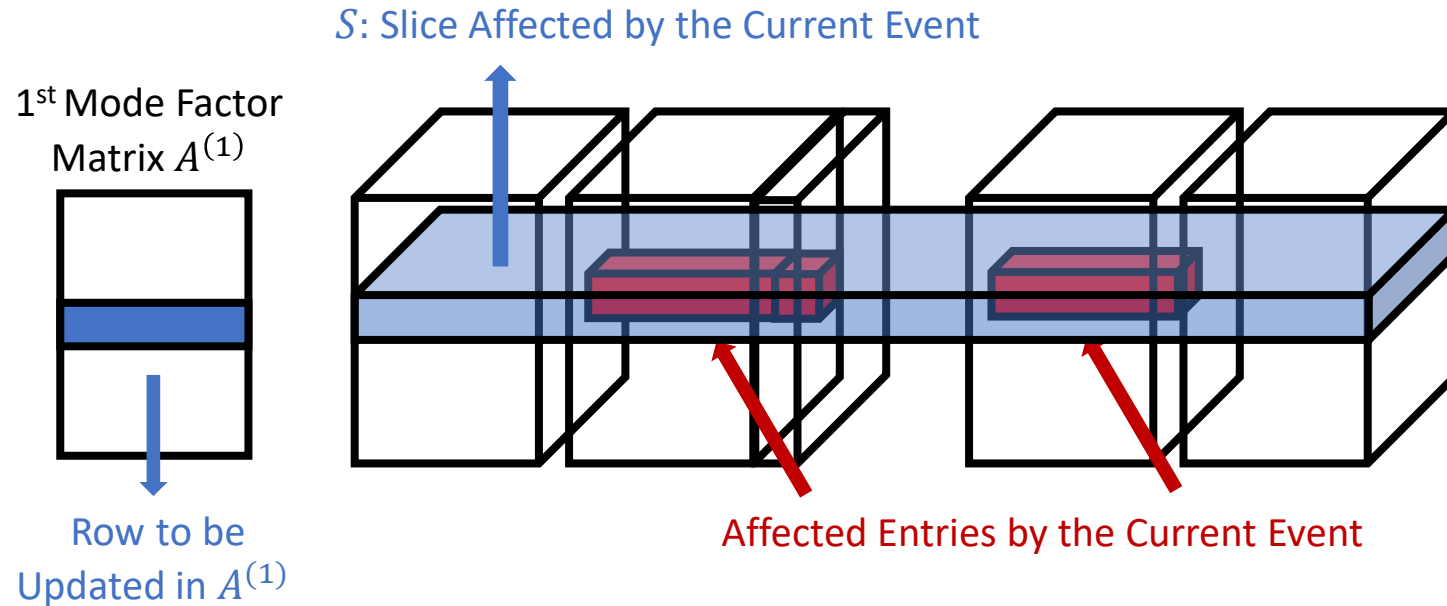
- **Approximate** \mathcal{X} as $\tilde{\mathcal{X}}$ (the **output** of **CP decomposition**) and solve the least square problem
- Computation is proportional to the number of non-zeros in $\Delta\mathcal{X}$



SliceNStitch-Vector: SNS_{VEC}

For non-time modes,

- Solve the original problem for only a single row
- Time complexity is proportional to the number of non-zeros in S



Pros

1. Significantly faster than SNS_{MAT}



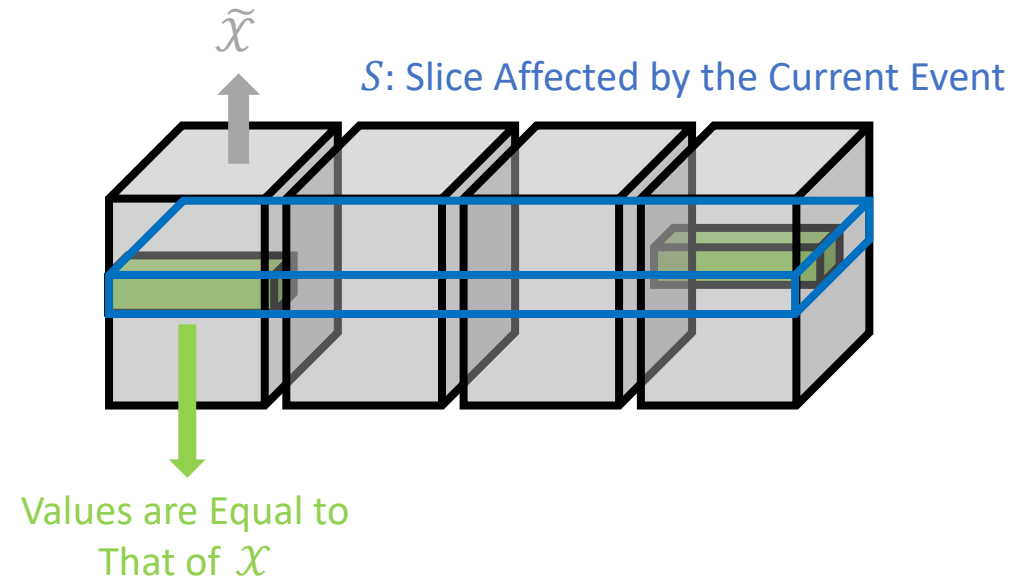
Cons

1. Numerically unstable
2. Slow downs if many non-zeros are of the same index



SliceNStitch-Random: SNS_{RND}

- Given a threshold θ
 - If the number of non-zeros in $S \leq \theta$ then solve the original problem
 - Else
 - Approximate \mathcal{X} as $\tilde{\mathcal{X}}$
 - correcting at most θ randomly chosen entries in $\tilde{\mathcal{X}}$ to \mathcal{X}
 - Solve the least square problem
- Time complexity is proportional to θ



Pros

1. Time complexity becomes constant



Cons

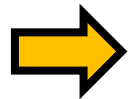
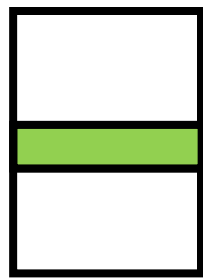
1. Numerically unstable
2. Reduction in the quality of the solution compared to SNS_{VEC}



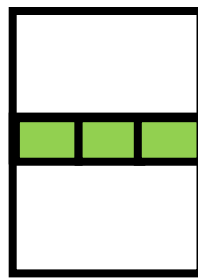
SliceNStitch-Stable: SNS_{VEC}^+ and SNS_{RND}^+

- Problem: SNS_{VEC} and SNS_{RND} are unstable
 - Many products (e.g. $\odot_{i \neq m} A^{(i)}$ and $*_{i \neq m} A^{(m)T} A^{(m)}$) are required
 - These result in too large numbers and impair the accuracy of the calculation
- Solution: update entries **one by one**, and **clip** each update value if it is larger than a threshold η

SNS_{VEC} and SNS_{RND}
: Update **all the entries**
in the row **at once**



SNS_{VEC}^+ and SNS_{RND}^+
: Update row entries **one**
by one



Guarantees that the updated entry does not increase the objective function although it is clipped

Road Map

- Introduction
- Problem Definition
- Data Model: Continuous Tensor Model
- Optimization Algorithms: SliceNStitch
- **Experiment Results** ←
- Conclusions



Experiments Settings

- 4 tensors from traffic and crime data

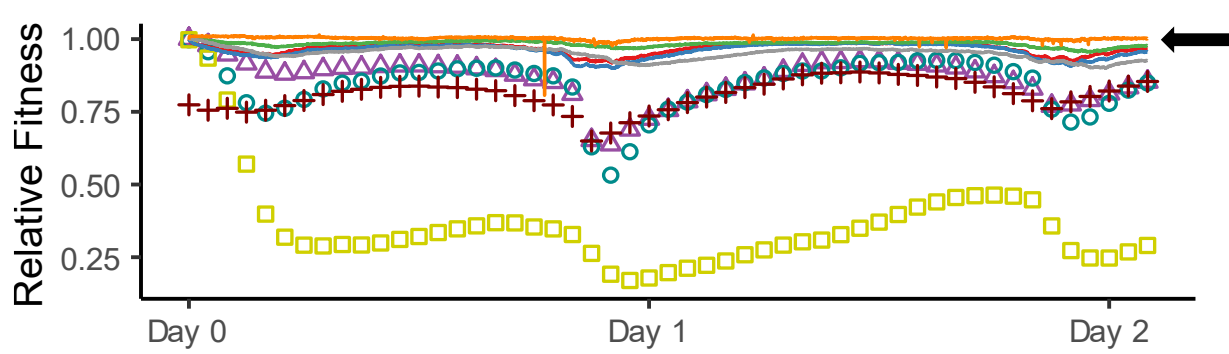
| Name | Size | # Non-zeros |
|---------------|-----------------------------------|-------------|
| Divvy Bikes | 673 x 673 x 525594 [minutes] | 3.82M |
| Chicago Crime | 77 x 32 x 148464 [hours] | 5.33M |
| New York Taxi | 265 x 265 x 5184000 [seconds] | 84.39M |
| Ride Austin | 219 x 219 x 24 x 285136 [minutes] | 0.89M |

- 4 baselines that update CPD periodically
 - ALS [CC70], onlineSCP [ZEB18], CP-stream [SHSK18], NeCPD [ASZ20]

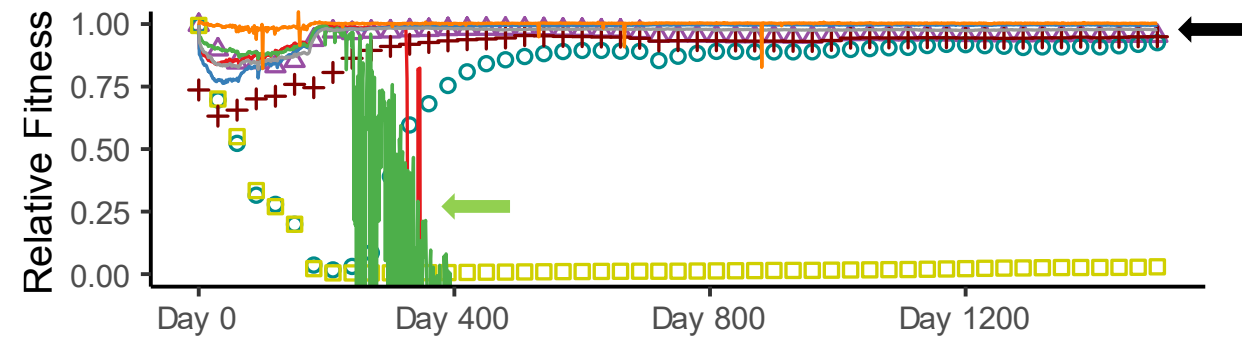
SliceNStitch is Accurate



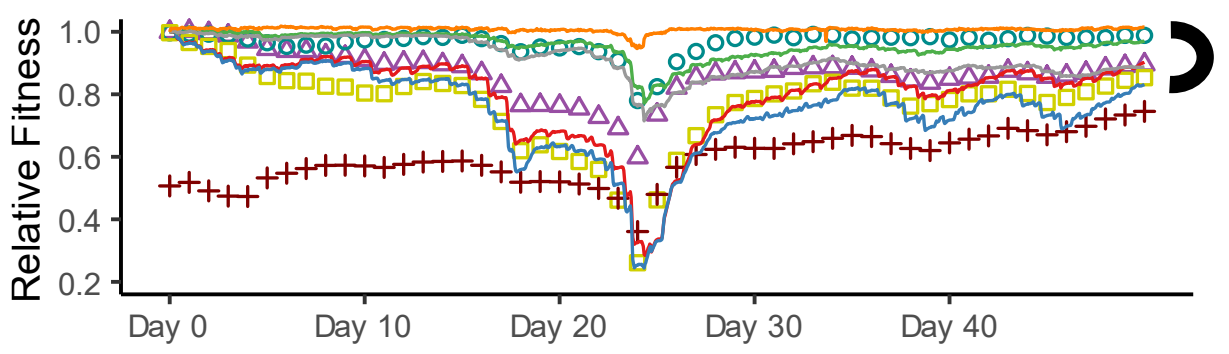
$$\text{Relative Fitness} = \frac{\text{Fitness}_{\text{target}}}{\text{Fitness}_{\text{ALS}}}$$



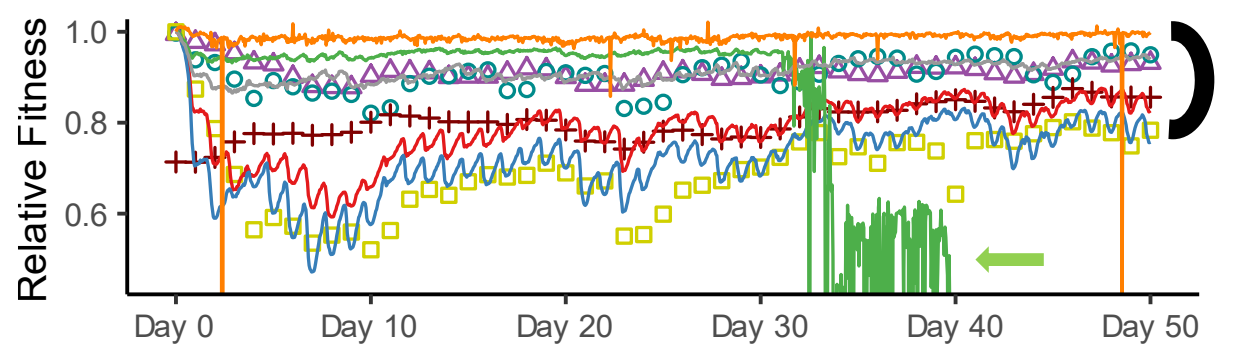
New York Taxi



Chicago Crimes

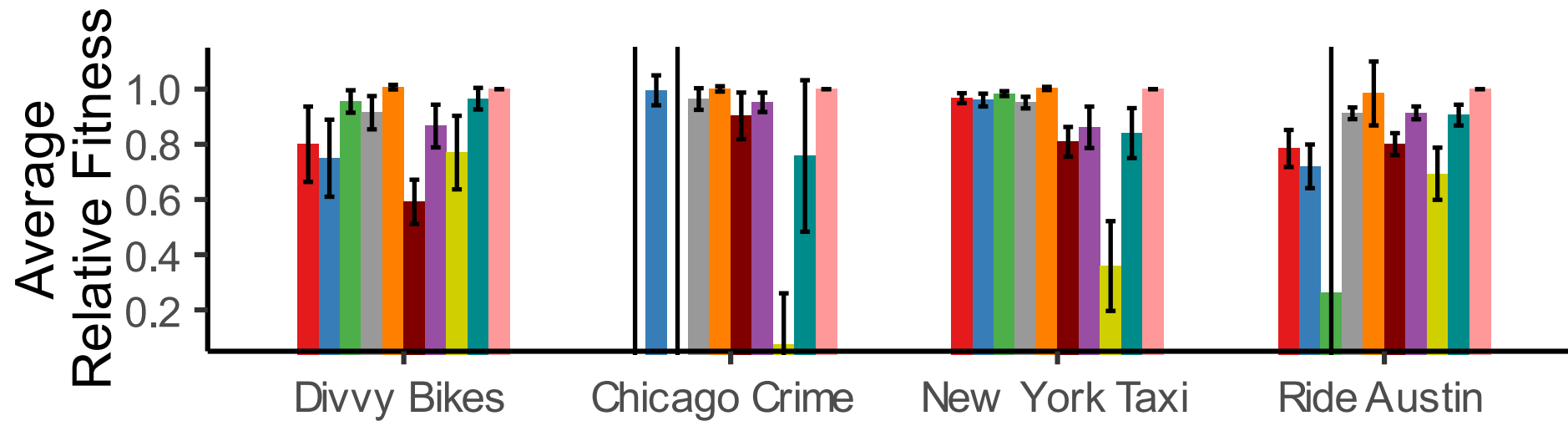


Divvy Bikes



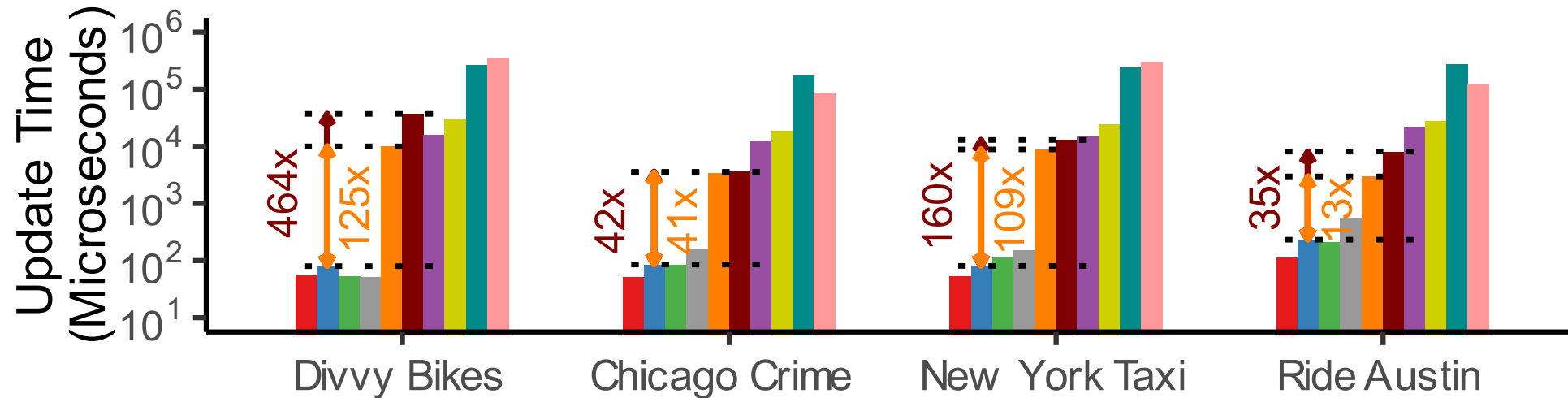
Ride Austin

SliceNStitch is Accurate



72 ~ 100% relative fitness to the most accurate baseline

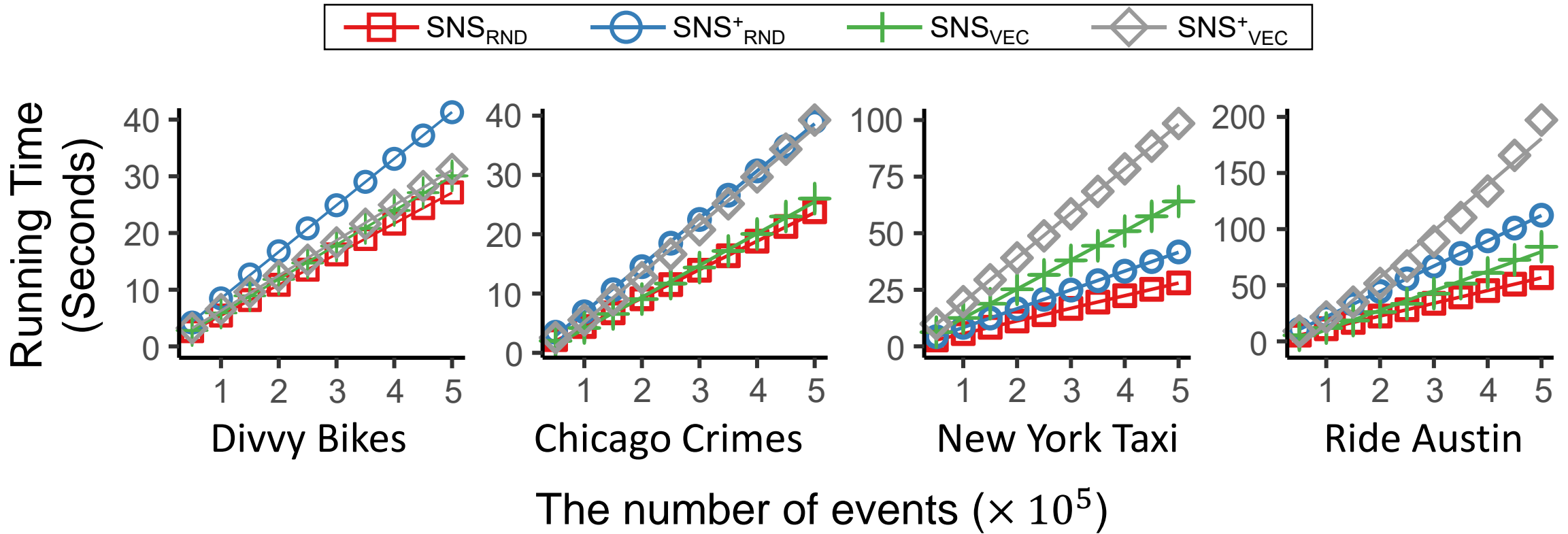
SliceNStitch is Fast



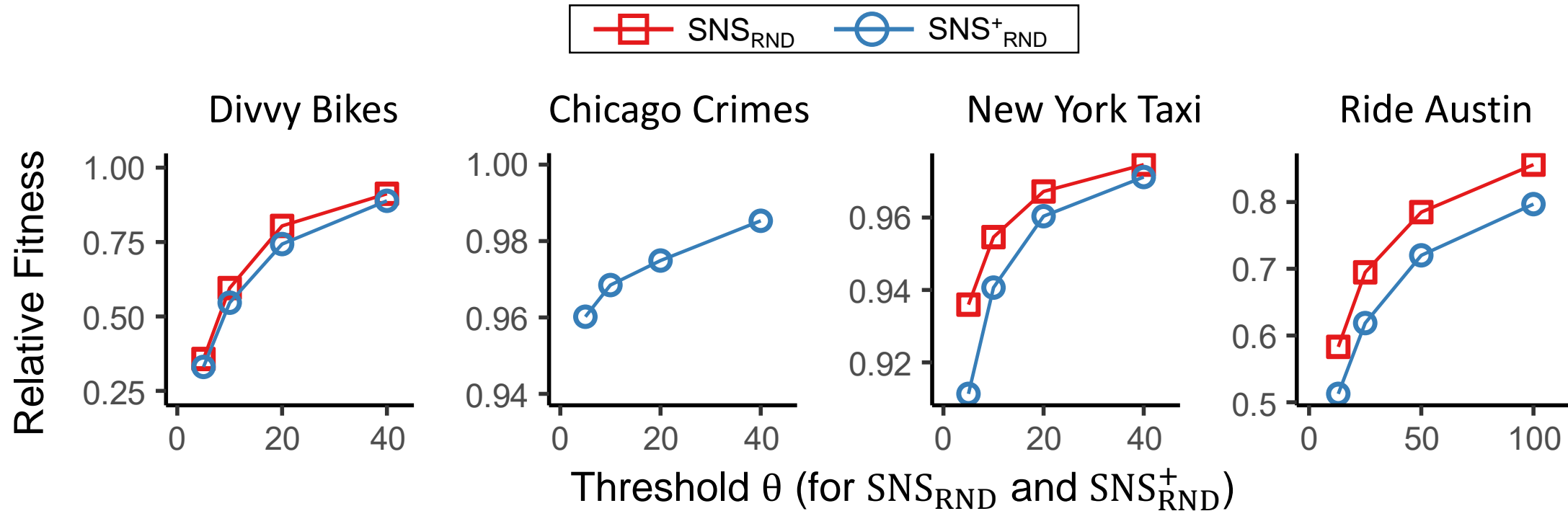
SNS^+_{RND} is up to 464 times faster than CP-stream

SliceNStitch is Scalable

The total runtime of all SliceNStitch versions was linear in the number of events.

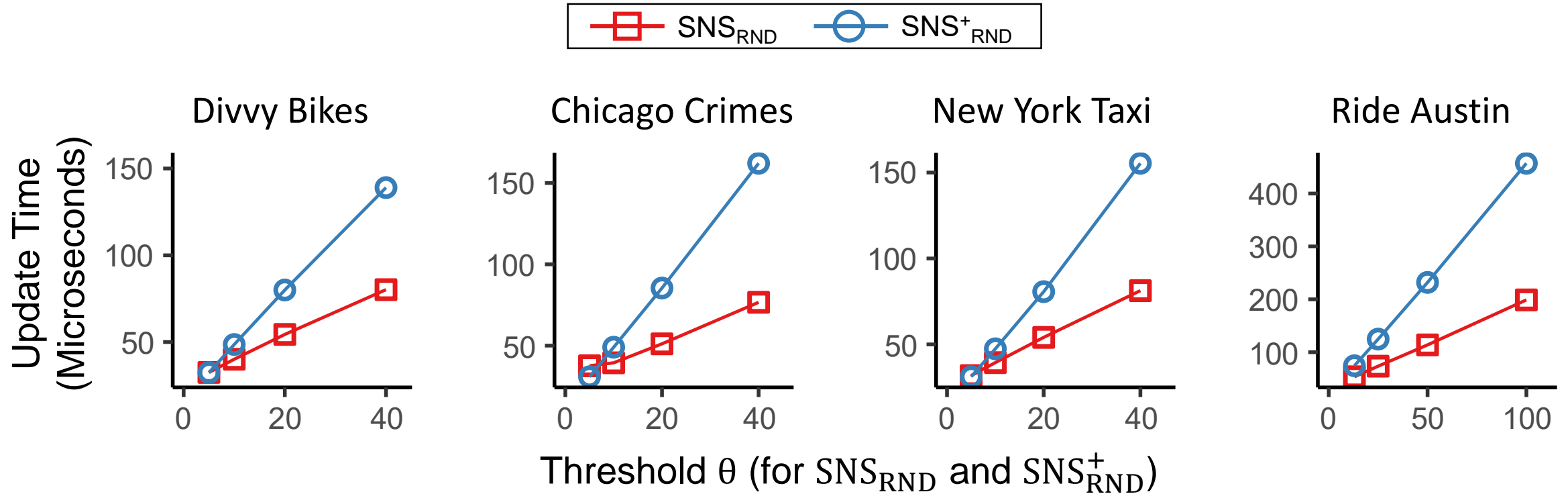


Effect of Sampling Parameter θ (for SNS_{RND} and SNS_{RND}^+)



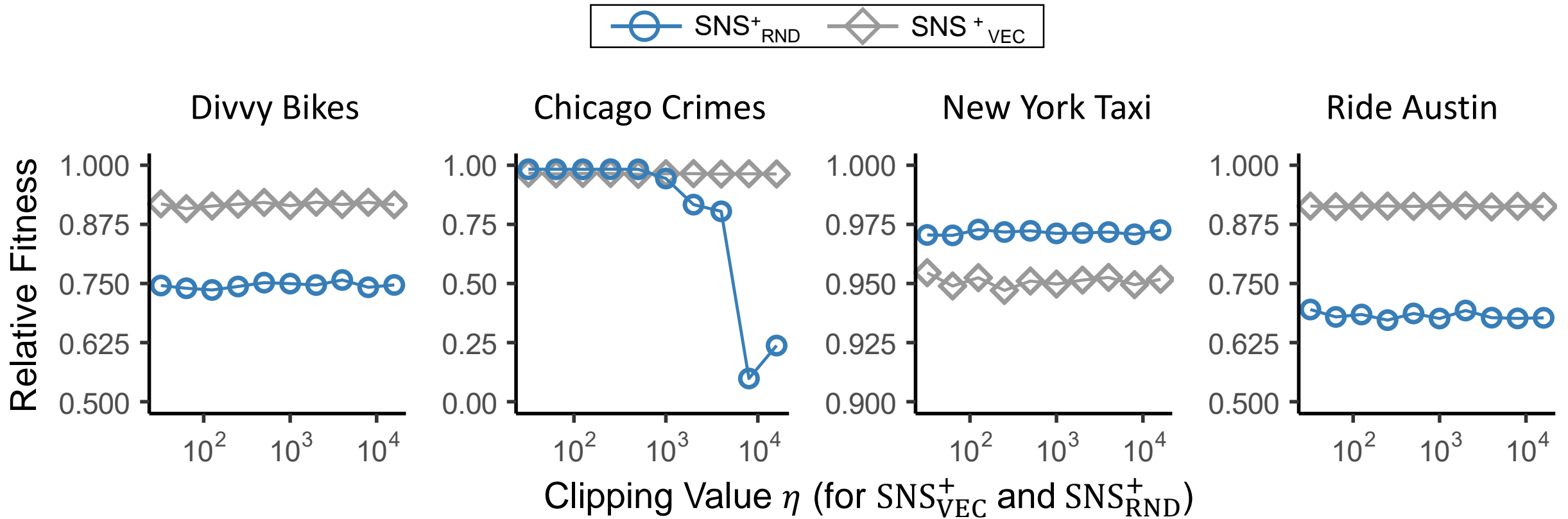
As θ increases, the fitness increases with diminishing returns

Effect of Sampling Parameter θ (for SNS_{RND} and SNS_{RND}^+)



As θ increases, runtime grows linearly

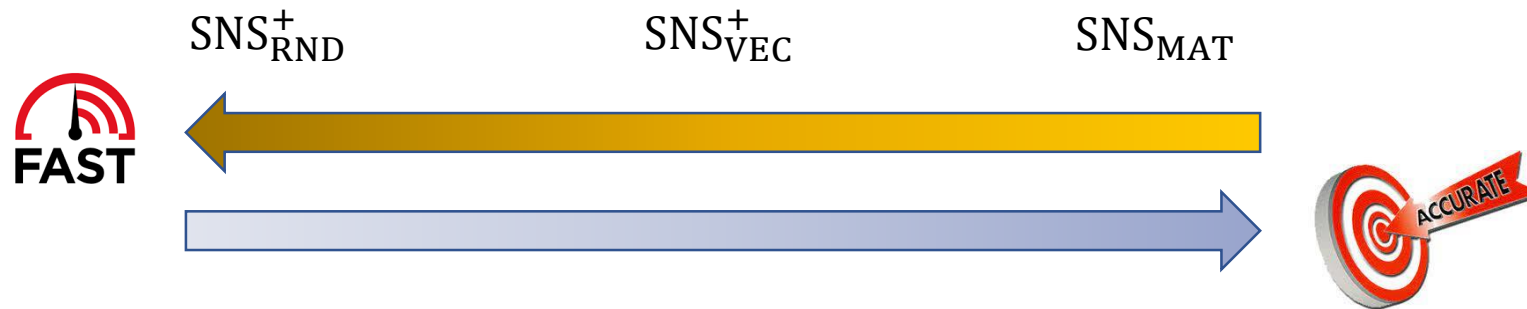
Effect of Clipping Value η (for $\text{SNS}_{\text{VEC}}^+$ and $\text{SNS}_{\text{RND}}^+$)



The fitness is insensitive to η as long as η is small enough

Practitioner's Guide

- We do not recommend SNS_{VEC} and SNS_{RND} due to numerical errors
- We recommend using the most accurate version within your runtime budget



- If SNS_{RND}^+ is chosen, increase θ enough within your runtime budget

SAMPLE

Road Map

- Introduction
- Problem Definition
- Data Model: Continuous Tensor Model
- Optimization Algorithms: SliceNStitch
- Experiment Results
- **Conclusions** ←



Conclusion

Continuous CPD with SliceNStitch achieved **near-instant updates, high fitness, and a small number of parameters.**

