

Data In, Facts Out:

Automated Monitoring of Facts by FactWatcher

Chengkai Li

Associate Professor, Department of Computer Science and Engineering

Director, Innovative Database and Information Systems Research (IDIR) Laboratory

University of Texas at Arlington

Nanjing University, Oct. 12th, 2015



Our Computational Journalism Project

Started in 2010. Collaborative project with Duke, Google Research, HP Labs, Stanford, and Chinese Academy of Sciences

- **Story finding:** finding and monitoring number-based facts pertinent to real-world events. The facts are leads to news stories.

FactWatcher

- **Fact checking:** discovering and checking factual claims in political discourses, social media, and news.

ClaimBuster



FactWatcher

Automated Monitoring of Facts from Real-
World Events



FactWatcher



Tuple t for new real world event appended to database



id	player	day	month	season	team	opp_team	pts	ast	reb
t_1	Bogues	11	Feb.	1991-92	Hornets	Hawks	4	12	5
t_2	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
t_3	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
t_4	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
t_5	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
t_6	Strictland	3	Jan.	1995-96	Blazers	Celtics	27	18	8
t_7	Wesley	25	Feb.	1995-96	Celtics	Nets	12	13	5

Find constraint-measure pair (C, M) such that t is in the contextual skyline



Generate factual claim



Wesley had 12 points, 13 assists and 5 rebounds on February 25, 1996 to become the first player with a 12/13/5 (points/assists/rebounds) in February.

Constraint	Measure
month=Feb	pts, ast, reb
opp_team=Nets	ast, reb
team=Celtics & opp_team=Nets	ast, reb
...	...

Fact Finding

Prominent streaks

Long consecutive subsequence of high values in a sequence

One-of-the-few objects

Qualifying statements that can only be made for very few objects

Situational facts

Comparison contexts and spaces that make a given object stand out



FactWatcher Finds Three Types of Facts (and can be Extended)

Domains

- sports, weather, crimes, transportation, finance, social media analytics

Examples from Real News Media

Prominent streaks

- “This month the Chinese capital has experienced 10 days with a maximum temperature in around 35 degrees Celsius – the most for the month of July in a decade.”

http://www.chinadaily.com.cn/china/2010-07/27/content_11055675.htm

- “The Nikkei 225 closed below 10000 for the 12th consecutive week, the longest such streak since June 2009.”

<http://www.bloomberg.com/news/articles/2010-08-06/japanese-stocks-fall-for-second-day-this-week-on-u-s-jobless-claims-yen>



FactWatcher Finds Three Types of Facts (and can be Extended)

Examples from Real News Media

Situational facts, One-of-the-few objects

- “Paul George had 21 points, 11 rebounds and 5 assists to become the first Pacers player with a 20/10/5 (points/rebounds/assists) game against the Bulls since Detlef Schrempf in December 1992.”
<http://espn.go.com/espn/elias?date=20130205>
- “The social world’s most viral photo ever generated 3.5 million likes, 170,000 comments and 460,000 shares by Wednesday afternoon.”
<http://www.cnbc.com/id/49728455>



FactWatcher Demo

<http://idir.uta.edu/factwatcher/>



»LIVE UPDATE

[February 20, 1998] **Todd Fuller** had 1 assist, 3 steals and 1 block in the Golden State Warriors' defeat against the Denver Nuggets. It is one of the best performance made by him.

SEARCH

michael jordan

Michael Adonis Jordan

Michael Jordan

Michael Michael Jordan

Michael Reggie Jordan

Michael Thomas Jordan

[January 13, 1997] **Horace Grant** had 26 points and 6 assists in the Orlando Magic's victory against the New Jersey Nets. It is one of the best performance made by him.

MORE LIKE THIS

[January 13, 1997] After the Orlando Magic's win over the New Jersey Nets, for the first time in his career, **Rony Seikaly** had at least 20 points for 6 consecutive games, after today's game.

MORE LIKE THIS

[January 13, 1997] **Horace Grant** had 26 points and 2 steals in the Orlando Magic's victory against the New Jersey Nets. It is one of the best performance made by him.

MORE LIKE THIS

[January 13, 1997] **Horace Grant** had 26 points, 6 assists and 2 steals in the Orlando Magic's victory against the New Jersey Nets. It is one of the best performance made by him.

MORE LIKE THIS

[January 13, 1997] After the Orlando Magic's victory against the New Jersey Nets, for the first time in his career, **Rony Seikaly** had at least 20 points and 8 rebounds for 6 consecutive games, after today's game.

MORE LIKE THIS

[January 13, 1997] **Nick Anderson** had 8 assists and 2 blocks in the Orlando Magic's win over the New Jersey Nets. It is one of the best performance made by him.

MORE LIKE THIS

FACT TYPE

SITUATIONAL FACT

PROMINENT STREAK

ONE-OF-THE-FEW

RANKING

RECENTNESS

INTERESTINGNESS

POPULARITY

PLAYERS

TEAMS

SEASONS

1996-97 (9)

1994-95 (5)

1992-93 (1)

+MORE

LESS-

Presented In



Excellent Demo Award

**COMPUTATION
+ JOURNALISM** 2014
SYMPOSIUM<http://idir.uta.edu/factwatcher/>

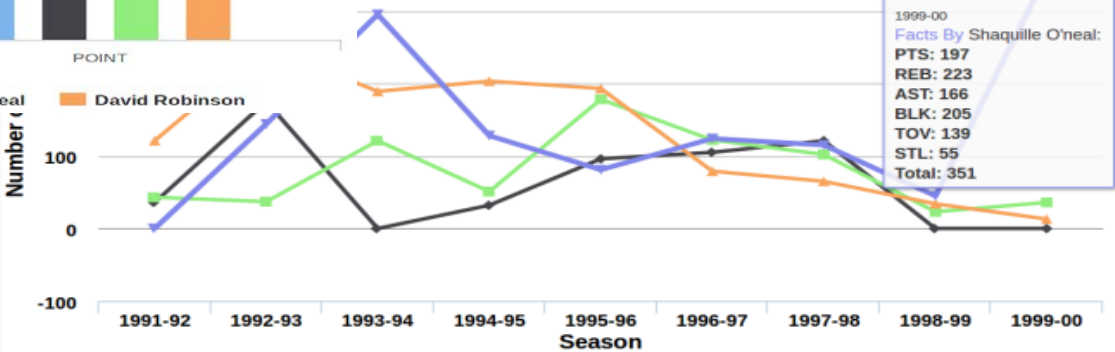
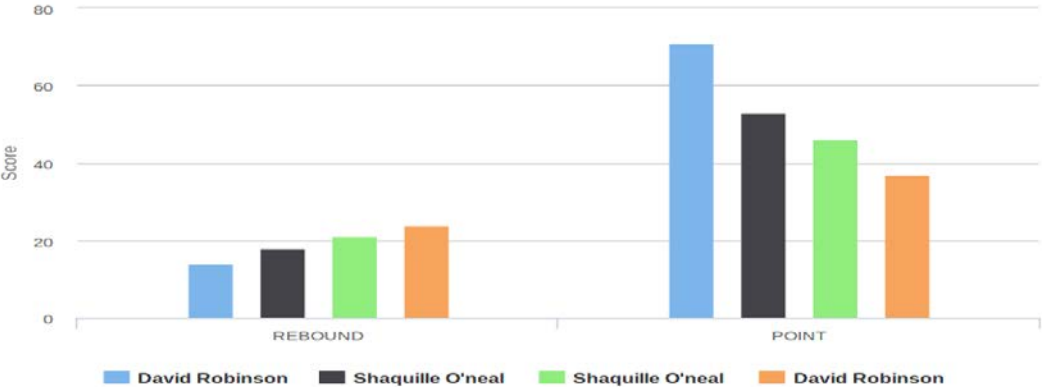
[April 24, 1994] **David Robinson** had 71 points and 14 rebounds in the San Antonio Spurs' victory against the Los Angeles Clippers. No one before had a better performance in NBA history.

[April 20, 1994] **Shaquille O'neal** had 53 points and 18 rebounds in the Orlando Magic's win over the Minnesota Timberwolves. No one before had a better performance in NBA history.

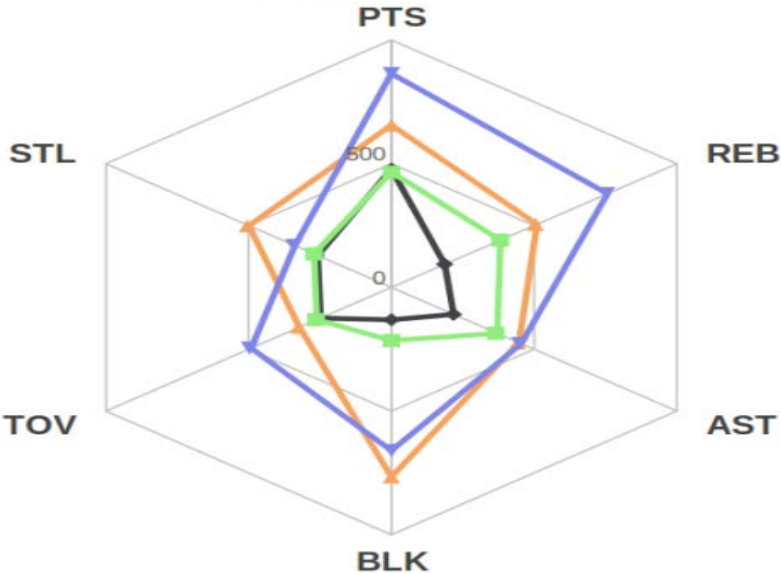
[February 16, 1993] **Shaquille O'neal** had 46 points and 21 rebounds in the Orlando Magic's defeat against the Detroit Pistons. No one before had a better performance in NBA history.

[February 27, 1992] **David Robinson** had 37 points and 24 rebounds in the San Antonio Spurs' victory against the Golden State Warriors. No one before had a better performance in NBA history.

Compare Similar Stories



Number of Facts



1999-00
Facts By Shaquille O'neal:
PTS: 197
REB: 223
AST: 166
BLK: 205
TOV: 139
STL: 55
Total: 351

How were these Facts Discovered in Current Systems?

Our (educated?) guess

- Experts monitor real-world events (e.g., watching an NBA game), have a gut-feeling, issue database queries, check out or not
- Prepared facts-to-be (e.g., Nowitzki only needs 477 more points to surpass O'Neal. Perhaps will happen around Christmas 2015)
- Predefined templates of facts/database queries
- Perhaps in-house systems/algorithms similar to FactWatcher



Elias Sports Bureau



StatSheet

No. 1-Seeded Louisville Cl

thevilledaily.com/louisville-basketball/game-recap/no-1-seeded-louisville-clips-no-4-seeded-michigan-82-76-wins-ncaa-championship

No. 1-Seeded Louisville Clips No. 4-Seeded Michigan 82-76, Wins NCAA Championship

Filed under [Game Recap](#) on April 9th, 2013

Share this recap



Tweet

Or

Like

One person likes this. Be the first of your friends.

NCAA Tournament 7th Round

	1ST	2ND	TOTAL	SPREAD	
 #4 Michigan	38	38	76	+4.0	●
 #1 Louisville	37	45	82	-4.0	●

Mon, Apr 08 2013, 10:23 PM EDT

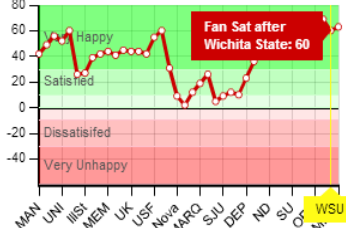
Georgia Dome
 Atlanta, Georgia
 Attendance: 74,326
 TV: CBS

[Boxscore](#) | [Game Notes](#) | [Game Recap](#) | [StatSmack](#)

No. 1-seeded Louisville got the win against No. 4-seeded Michigan 82-76 in the Championship Game of the NCAA Tournament on Monday, Apr. 8. The Cardinals were led by Peyton Siva, who got 18 points and six rebounds (5 Ast 4 Stl). Gorgui Dieng also had an outstanding outing, scoring eight points and adding eight rebounds (6 Ast 3 Blk). Michigan closes out its impressive season with a 31-8 overall record. The Wolverines got to the NCAA Tournament as an at-large team after falling to Wisconsin 68-59 in the Big Ten Tournament. In the regular season, they finished fourth in the Big Ten with a 12-6 conference record. In making the national championship game, Michigan knocked off No. 13-seeded South Dakota State 71-56 in the second round and No. 5-seeded Virginia Commonwealth 78-53 in the third round. Following that, the Wolverines got through No. 1-seeded Kansas 87-85 in the Sweet Sixteen, No. 3-seeded Florida 79-59 in the Elite Eight, and No. 4-seeded Syracuse 61-56 in the Final Four. For the Wolverines, Trey Burke got a game-high 24 points and four rebounds. Michigan (31-8) finished the regular season fourth in the Big Ten with a 12-6 record. Through their amazing run, Louisville got through No. 16-seeded North Carolina A&T 79-48 in the second round and No. 8-seeded Colorado State 82-56 in the third round. Following that, the Cardinals got through No. 12-seeded Oregon 77-69 in the Sweet Sixteen, No. 2-seeded Duke 85-63 in the Elite Eight, and No. 9-seeded Wichita State 72-68 in the Final Four.


StatSeed: NCAA Automatic #1 Seed

Fan Satisfaction



More about Fan Satisfaction

Find another NCAA team:



Categories

©2015 The University of Texas at Arlington. All Rights Reserved.

Narrative Science

Forbes Earnings Preview: / X

www.forbes.com/sites/narrativescience/2013/05/02/forbes-earnings-preview-anadarko-petroleum-6/

Forbes

New Posts
+18 posts this hour

Popular
Most Reputable Comp

Lists
The Global 2000

Video
Hip-Hop's Wealthiest

Search

Narrative Science
Forbes Partner
+ Follow (50)

xerox
Ready For Real Business

INVESTING | 5/02/2013 @ 2:31PM | 488 views

Forbes Earnings Preview:
Anadarko Petroleum

By Narrative Science

+ Comment Now + Follow Comments

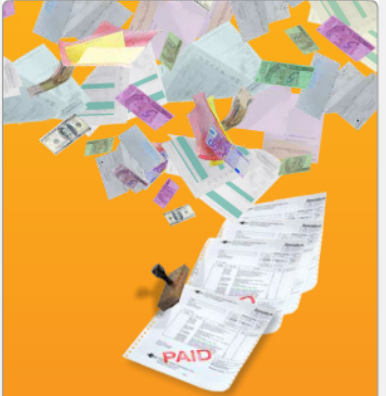
Analysts have become increasingly bullish on [Anadarko Petroleum](#) APC +2.02% (APC) in the month leading up to the company's first quarter earnings announcement scheduled for Monday, May 6, 2013. The consensus earnings per share estimate has moved up from 88 cents a share to the current expectation of earnings of 91 cents a share.

[Wall Street](#) projections are down 1.1% year-over-year, as the company reported earnings of 92 cents per share.

The consensus estimate has gone up, from 82 cents, over the past three months. Analysts are expecting earnings of \$4.04 per share for the fiscal year. Revenue is projected to be \$3.49 billion for the quarter, 1.2% above the year-earlier total of \$3.45 billion. For the year, revenue is projected to roll in at \$15.21 billion.

Revenue has declined for the third quarter in a row. The year-over-

0
Share
13
Tweet
1
Share
0
Submit
0
+1



Handling \$421 billion in accounts payables annually for companies like yours.

[learn more ▶](#)

Ready For Real Business xerox

©2015 The University of Texas at Arlington. All Rights Reserved.

Publications

- [Online Frequent Episode Mining](#). Xiang Ao, Ping Luo, Chengkai Li, Fuzhen Zhuang, and Qing He. ICDE 2015, pages 891-902.
- [Data In, Fact Out: Automated Monitoring of Facts by FactWatcher](#). Naeemul Hassan, Afroza Sultana, You Wu, Gensheng Zhang, Chengkai Li, Jun Yang, and Cong Yu. VLDB 2014, pages 1557-1560. Demonstration description. (**excellent demonstration award**)
- [Finding, Monitoring, and Checking Claims Computationally Based on Structured Data](#). Brett Walenz, You (Will) Wu, Seokhyun (Alex) Song, Emre Sonmez, Eric Wu, Kevin Wu, Pankaj K. Agarwal, Jun Yang, Naeemul Hassan, Afroza Sultana, Gensheng Zhang, Chengkai Li, Cong Yu. 2014 Computation+Journalism Symposium.
- [Incremental Discovery of Prominent Situational Facts](#). Afroza Sultana, Naeemul Hassan, Chengkai Li, Jun Yang, Cong Yu. ICDE 2014, pages 112-123.
- [Discovering General Prominent Streaks in Sequence Data](#). Gensheng Zhang, Xiao Jiang, Ping Luo, Min Wang, Chengkai Li. ACM TKDD, 8(2):article 9, June 2014.
- [Discovering and Learning Sensational Episodes of News Events](#). Xiang Ao, Ping Luo, Chengkai Li, Fuzhen Zhuang, Qing He, and Zhongzhi Shi. WWW 2014, pages 217-218.
- [On "One of the Few" Objects](#). You Wu, Pankaj K. Agarwal, Chengkai Li, Jun Yang, Cong Yu. KDD 2012, pages 1487-1495.
- [Prominent Streak Discovery in Sequence Data](#). Xiao Jiang, Chengkai Li, Ping Luo, Min Wang, Yong Yu. KDD 2011, pages 1280-1288.



Incremental Discovery of Prominent Situational Facts. Afroza Sultana, Naeemul Hassan, Chengkai Li, Jun Yang, Cong Yu. ICDE 2014, pages 112-123.



Situational Facts

“Paul George had 21 points, 11 rebounds and 5 assists to become the first Pacers player with a 20/10/5 (points/rebounds/assists) game against the Bulls since Detlef Schrempf in December 1992.”

(<http://espn.go.com/espn/elias?date=20130205>)



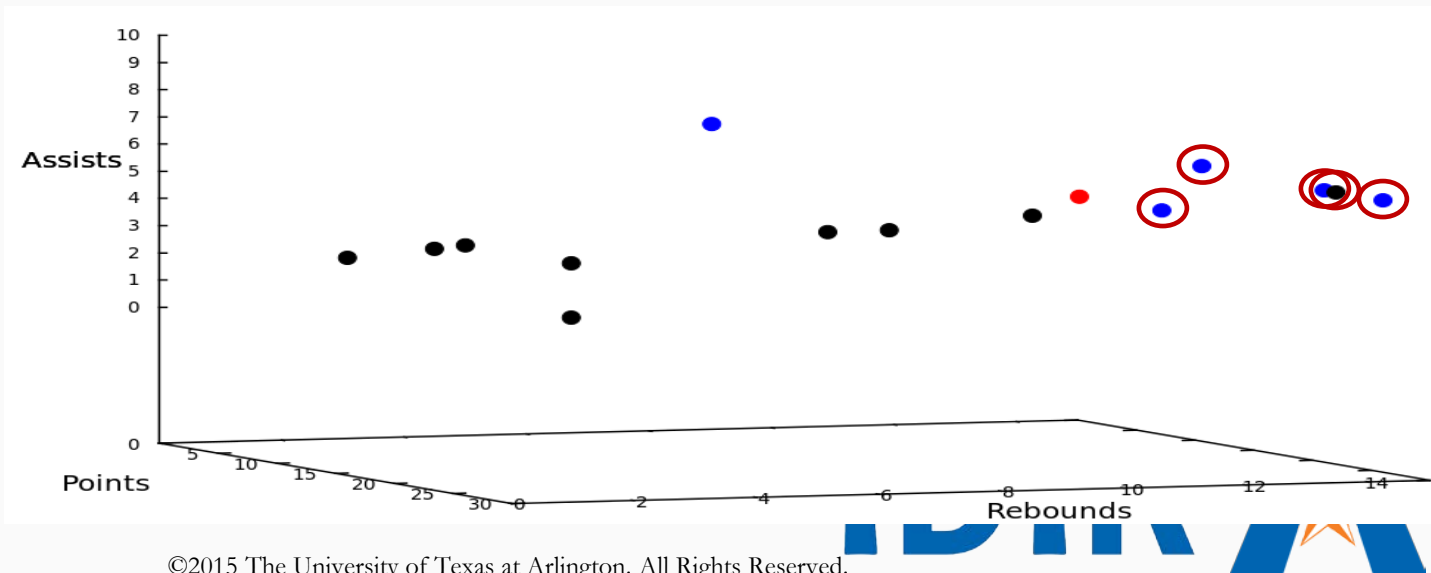
Skyline



Situational Facts

“Paul George had 21 points, 11 rebounds and 5 assists to become the first Pacers player with a 20/10/5 (points/rebounds/assists) game against the Bulls since Detlef Schrempf in December 1992.”

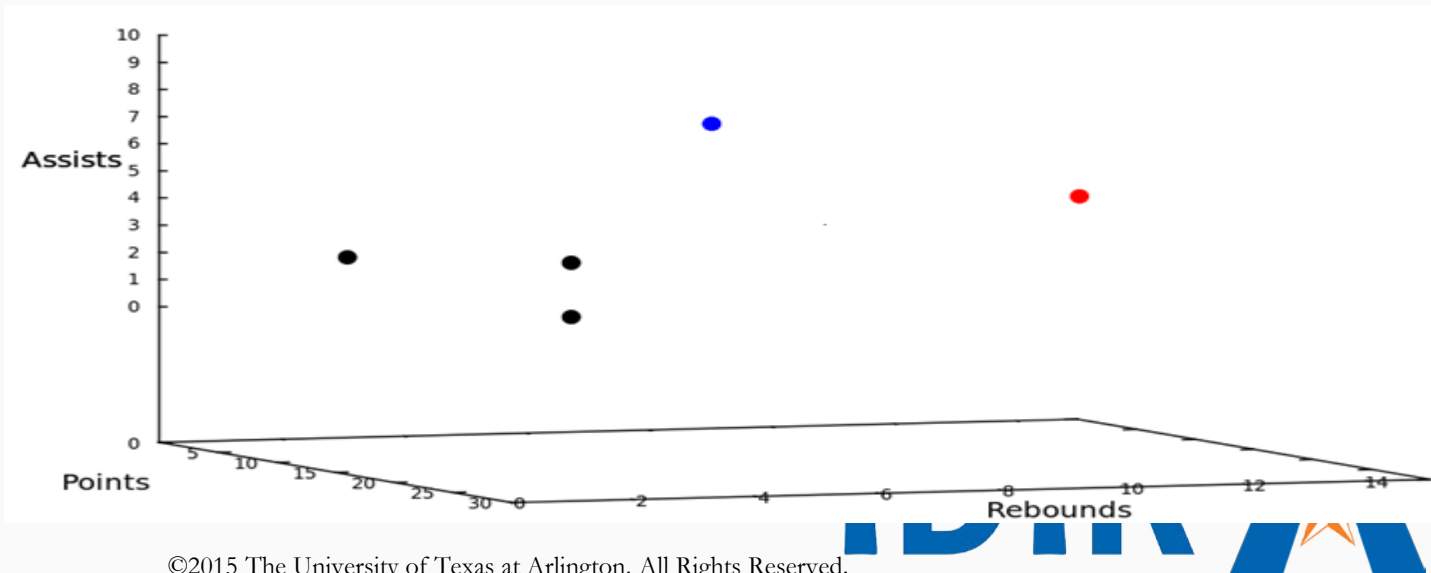
(<http://espn.go.com/espn/elias?date=20130205>)



Situational Facts

“Paul George had 21 points, 11 rebounds and 5 assists to become the first Pacers player with a 20/10/5 (points/rebounds/assists) game against the Bulls since Detlef Schrempf in December 1992.”

(<http://espn.go.com/espn/elias?date=20130205>)



Situational Facts

“The social world’s most viral photo ever generated 3.5 million likes, 170,000 comments and 460,000 shares by Wednesday afternoon.”

(<http://www.cnn.com/id/49728455/President Obama Sets New Social Media Record>)



Situational Facts

“The social world’s most viral photo ever generated 3.5 million likes, 170,000 comments and 460,000 shares by Wednesday afternoon.”

(<http://www.cnn.com/id/49728455/President Obama Sets New Social Media Record>)



Situational Facts

“The social world’s **most viral photo** ever generated **3.5 million likes, 170,000 comments and 460,000 shares** by Wednesday afternoon.”

(<http://www.cnn.com/id/49728455/President Obama Sets New Social Media Record>)



Situational Facts

- **Stock Data:** Stock A becomes the first stock in history with price over \$300 and market cap over \$400 billion.
- **Weather Data:** Today's measures of wind speed and humidity are x and y, respectively. City B has never encountered such high wind speed and humidity in March.
- **Criminal Records:** There were 50 DUI arrests and 20 collisions in city C yesterday, the first time in 2013.

Financial Analyst
Journalists
Scientists Citizens



A Mini-world of Basketball GameLogs

id	player	day	month	season	team	opp_team	pts	ast	reb
t_1	Bogues	11	Feb.	1991-92	Hornets	Hawks	4	12	5
t_2	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
t_3	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
t_4	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
t_5	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
t_6	Strictland	3	Jan.	1995-96	Blazers	Celtics	27	18	8
t_7	Wesley	25	Feb.	1995-96	Celtics	Nets	12	13	5

Last tuple appended to table



A Mini-world of Basketball GameLogs

id	player	day	month	season	team	opp_team	pts	ast	reb
t_1	Bogues	11	Feb.	1991-92	Hornets	Hawks	4	12	5
t_2	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
t_3	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
t_4	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
t_5	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
t_6	Strictland	3	Jan.	1995-96	Blazers	Celtics	27	18	8
t_7	Wesley	25	Feb.	1995-96	Celtics	Nets	12	13	5



A Mini-world of Basketball GameLogs

id	player	day	month	season	team	opp_team	pts	ast	reb
t_1	Bogues	11	Feb.	1991-92	Hornets	Pacers	4	12	5
t_2	Seakaly	13	Feb.	1991-92	Hornets	Pacers	24	5	15
	Sherman	7	Dec.	1993-94	Celtics	Nets			
t_4	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
t_5	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
	Stinetland	8	Jun.	1995-96	Blazers	Celtics			
t_7			Feb.				12	13	5

A Mini-world of Basketball GameLogs

id	player	day	month	season	team	opp_team	pts	ast	reb
t_1	Diggins	11	Feb.	1991-92	Hornets	Pacers	4	12	5
t_2	Sarkis	12	Feb.	1991-92	Hornets	Pacers	24	5	15
	Sherman			1993-94	Celtics	Pacers			
t_4	Wesley	25	Feb.	1996-97	Celtics	Nets	2	5	2
t_5	Wesley		Feb.	1996-97	Celtics	Timberwolves	3	5	3
	Stapleton			1997-98	Pacers	Celtics			
t_7			Feb.				12	13	5

■ Wesley had 12 points, 13 assists and 5 rebounds on February 25, 1996 to become the first player with a 12/13/5 (points/assists/rebounds) in February.



A Mini-world of Basketball GameLogs

id	player	day	month	season	team	opp_team	pts	ast	reb
t_1	Bogues	11	Feb	1991-92	Jazz	Knicks	1	1	5
t_2	Serkaly	13	Feb	1991-92	Jazz	Knicks	1	5	15
t_3	Sherman	7	Dec	1993-94	Knicks	Knicks	1	1	5
t_4	Wesley	4	Feb	1994-95	Knicks	Nets	1	5	2
t_5	Wesley	5	Feb	1994-95	Knicks	Timberwolves	1	5	1
t_6	Stinetland	9	Jan	1995-96	Pacers	Celtics	27	18	8
t_7				1995-96			12	13	5

A Mini-world of Basketball GameLogs

id	player	day	month	season	team	opp_team	pts	ast	reb
	Dumars	11	Feb	1991-92	Dumars	Dumars	11		
	Sarkis	15	Feb	1991-92	Dumars	Dumars	11		
t_3	Shannon		Dec	1993-94	Celtics	Nets		13	5
t_4	Wesley	25	Feb	1996-97	Celtics	Nets		5	2
	Wesley	25	Feb	1996-97	Celtics	Dumars			
	Shannon		Dec	1993-94	Dumars	Celtics			
t_7					Celtics	Nets		13	5

- Wesley had 13 assists and 5 rebounds on February 25, 1996 to become the second Celtics player with a 13/5 (assists/rebounds) game against the Nets.



Problem Definition

Dimension space: $\mathcal{D} = \{d_1, \dots, d_n\}$

Measure space: $\mathcal{M} = \{m_1, \dots, m_s\}$

id	player	day	month	season	team	opp_team	pts	ast	reb
t_1	Bogues	11	Feb.	1991-92	Hornets	Hawks	4	12	5
t_2	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
t_3	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
t_4	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
t_5	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
t_6	Strictland	3	Jan.	1995-96	Blazers	Celtics	27	18	8

append-only table



Problem Definition

□ **Constraint (C):** $d_1=v_1 \wedge d_2=v_2 \wedge \dots \wedge d_n=v_n, v_i \in \text{dom}(d_i) \cup \{*\}$

- $\text{team}=\text{Celtics} \wedge \text{opp_team}=\text{Nets}$

id	player	day	month	season	team	opp_team	pts	ast	reb
t_1	Bogues	11	Feb	1991-92	Knicks	Knicks	12	1	5
t_2	Seikaly	13	Feb	1991-92	Knicks	Knicks	11	5	15
t_3	Sherman		Dec	1993-94	Celtics	Nets	11	1	5
t_4	Westley		Feb	1994-95	Celtics	Nets			
t_5	Westley	5	Feb	1994-95	Celtics	Knicks	5		
t_6	Strickland	8	Jan	1995-96	Bruins	Celtics		15	8

Problem Definition

□ **Constraint-Measure Pair (C, M)** : Combination of a constraint and measure subspace

- $(\text{team}=\textit{Celtics} \wedge \text{opp_team}=\textit{Nets}, \{\text{assists}, \text{rebounds}\})$

id	player	day	month	season	team	opp_team	pts	ast	reb
	Bogues	11	Feb	1991-92	Bogues	Bogues	15		
	Sorkin	14	Feb	1991-92	Heat	Heat	15		
t_3	Sherman		Dec	1993-94	Celtics	Nets		13	5
t_4	Westley		Feb	1994-95	Celtics	Nets		5	2
	Westley	5	Feb	1994-95	Celtics	Knicks			
	Strickland		Jun	1995-96	Brazers	Celtics			

Problem Definition

□ **Contextual skyline:** skyline regarding (C, M)

- $\sigma_{\text{team}=\text{Celtics} \wedge \text{opp_team}=\text{Nets}}(R), M=\{\text{assists,rebounds}\}$
➤ $\{t_3\}$

id	player	day	month	season	team	opp_team	pts	ast	reb
t_1	Bogues	11	Feb	1991-92	Pacers	Pacers	15	1	5
t_2	Seikaly	13	Feb	1991-92	Pacers	Pacers	15	1	15
t_3	Sherman	7	Dec	1993-94	Celtics	Nets	15	13	5
t_4	Wesley	4	Feb	1994-95	Celtics	Nets	15	5	2
t_5	Wesley	5	Feb	1994-95	Celtics	Timberwolves	15	1	1
t_6	Strickland	8	Jan	1995-96	Pacers	Celtics	15	15	8

FactWatcher



Tuple t for new real world event appended to database



id	player	day	month	season	team	opp_team	pts	ast	reb
t_1	Bogues	11	Feb.	1991-92	Hornets	Hawks	4	12	5
t_2	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
t_3	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
t_4	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
t_5	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
t_6	Strictland	3	Jan.	1995-96	Blazers	Celtics	27	18	8
t_7	Wesley	25	Feb.	1995-96	Celtics	Nets	12	13	5

Find constraint-measure pair (C, M) such that t is in the contextual skyline



Generate factual claim



Wesley had 12 points, 13 assists and 5 rebounds on February 25, 1996 to become the first player with a 12/13/5 (points/assists/rebounds) in February.

Constraint	Measure
month=Feb	pts, ast, reb
opp_team=Nets	ast, reb
team=Celtics & opp_team=Nets	ast, reb
...	...

Related Work

- Conventional skyline analysis (Borzsonyi et al. ICDE 2001)
 - Q : context, measure subspace $\implies A$: contextual skyline tuples
 - ✓ Our focus--- A : tuple $\implies Q$: constraint-measure pairs



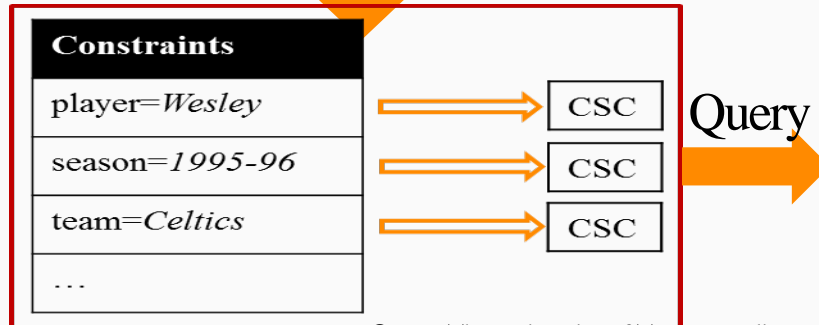
Related Works

➤ Compressed Skycube (Xia et al. SIGMOD 2006)

- Update compressed skycube in monitoring fashion

✓ We adapted CSC for each constraint: **Constraint-CSC**

id	player	day	month	season	team	opp_team	pts	ast	reb
t_1	Bogues	11	Feb.	1991-92	Hornets	Hawks	4	12	5
t_2	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
t_3	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
t_4	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
t_5	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
t_6	Strickland	3	Jan.	1995-96	Blazers	Celtics	27	18	8
t_7	Wesley	25	Feb.	1995-96	Celtics	Nets	12	13	5



Constraint	Measure
month=Feb	pts, ast, reb
opp_team=Nets	ast, reb
team=Celtics & opp_team=Nets	ast, reb
...	...

Related Works

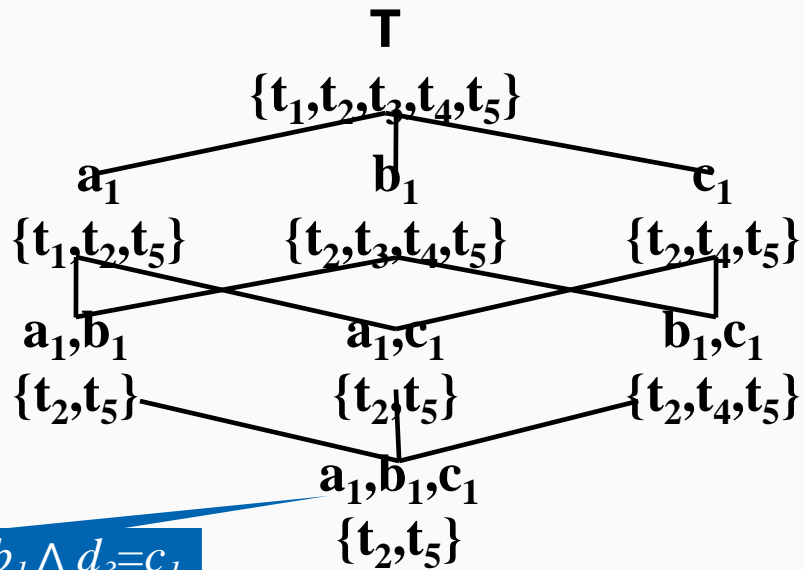
➤ Prominent Analysis by Ranking (Wu et. Al. VLDB 2009)

- Static data, onetime query
 - ✓ We dealt on continuous data, standing query
- Find the contexts where an object is ranked high in a **single scoring attribute**
 - ✓ We considered skyline on **multiple measure subspaces**



Modeling

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



$$d_1=a_1 \wedge d_2=b_1 \wedge d_3=c_1$$

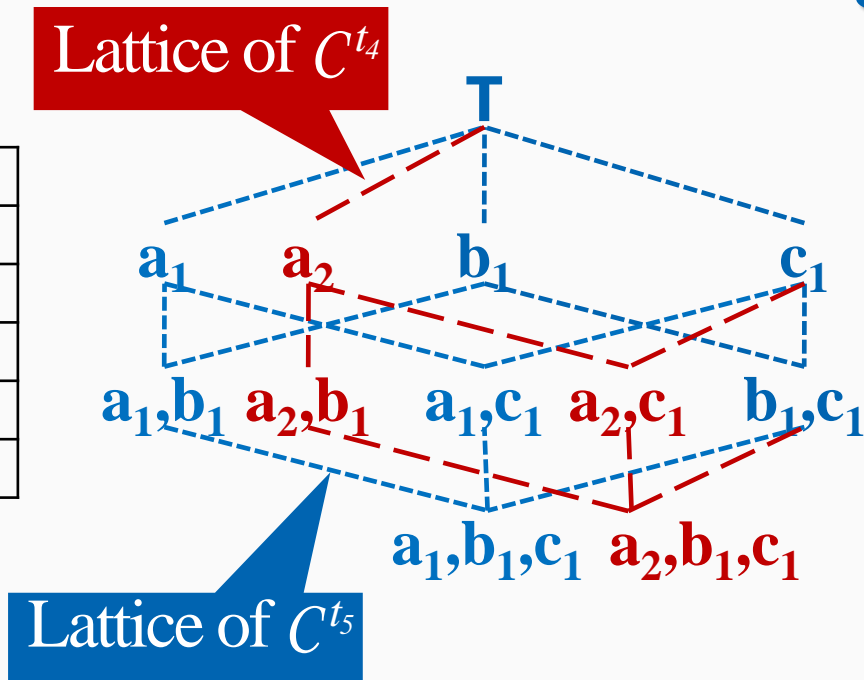
Lattice of C^t

Tuple Satisfied Constraint C^t : If $\forall d_i \in \mathcal{D}$, $C.d_i = *$ or $C.d_i = t.d_i$, t satisfies C .



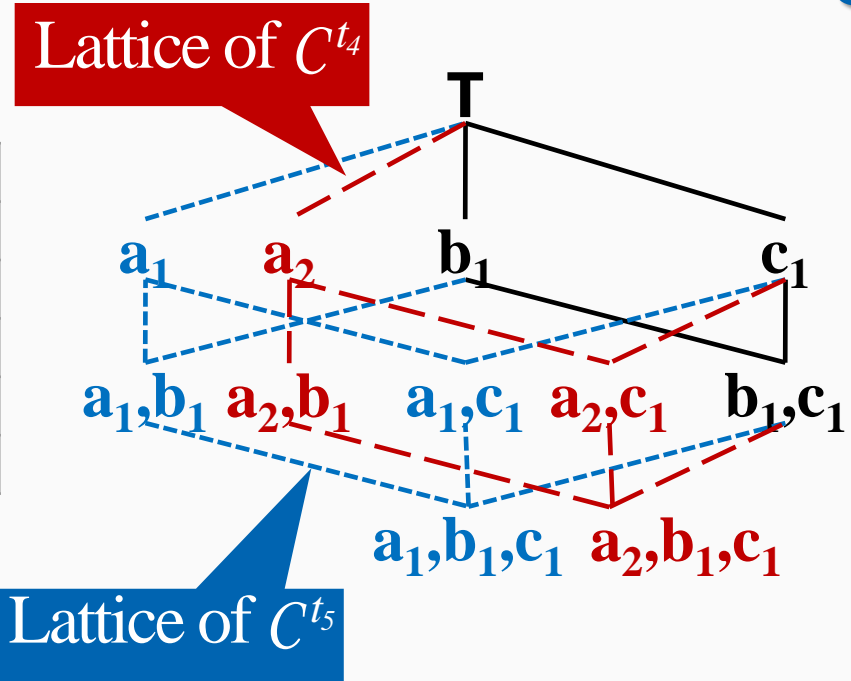
Modeling

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



Modeling

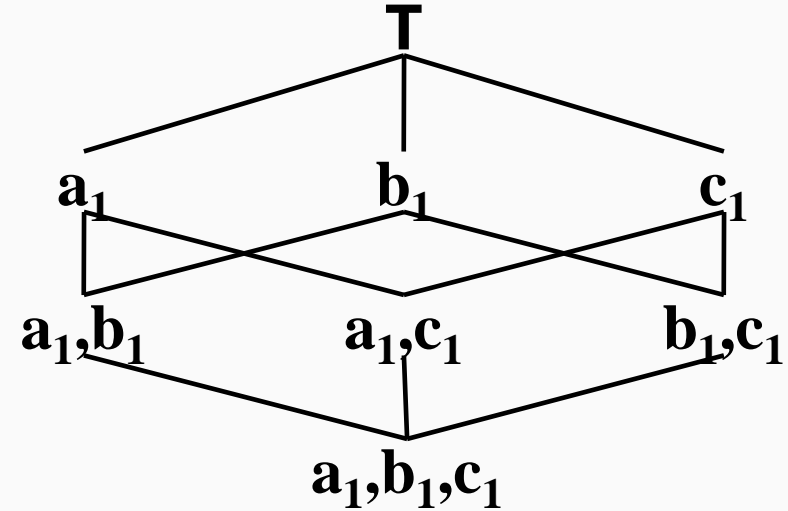
id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



Lattice Intersection: $C^{t_4 t_5} = C^{t_4} \cap C^{t_5}$

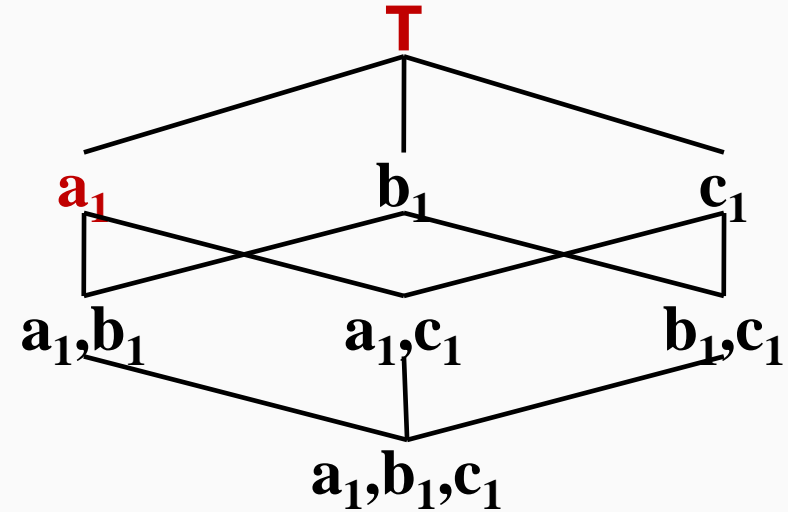
Brute-Force Approach

<i>id</i>	<i>d₁</i>	<i>d₂</i>	<i>d₃</i>	<i>m₁</i>	<i>m₂</i>
<i>t₁</i>	<i>a₁</i>	<i>b₂</i>	<i>c₂</i>	10	15
<i>t₂</i>	<i>a₁</i>	<i>b₁</i>	<i>c₁</i>	15	10
<i>t₃</i>	<i>a₂</i>	<i>b₁</i>	<i>c₂</i>	17	17
<i>t₄</i>	<i>a₂</i>	<i>b₁</i>	<i>c₁</i>	20	20
<i>t₅</i>	<i>a₁</i>	<i>b₁</i>	<i>c₁</i>	11	15



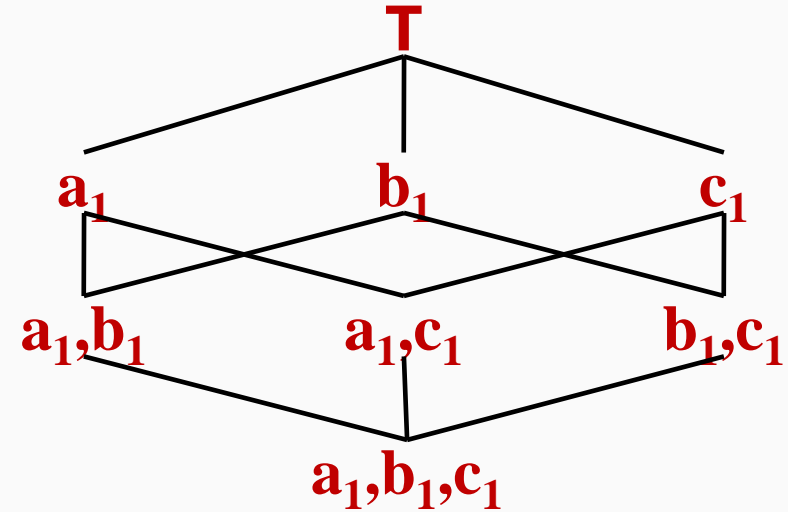
Brute-Force Approach

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



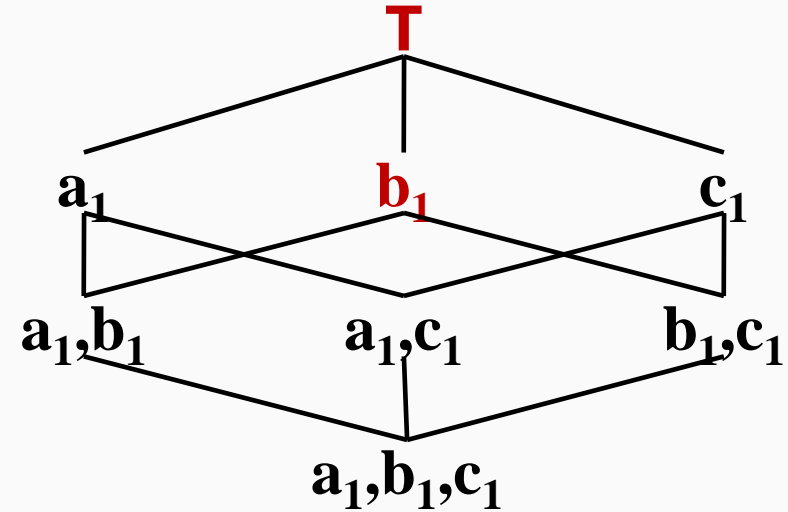
Brute-Force Approach

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



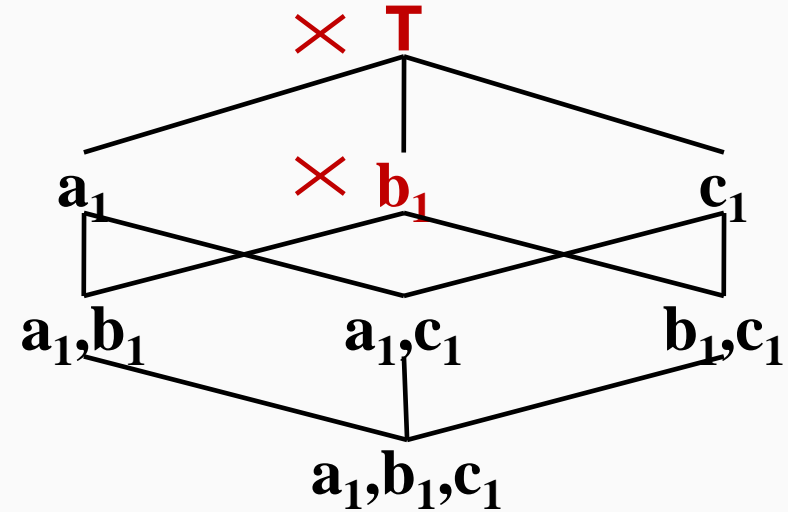
Brute-Force Approach

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



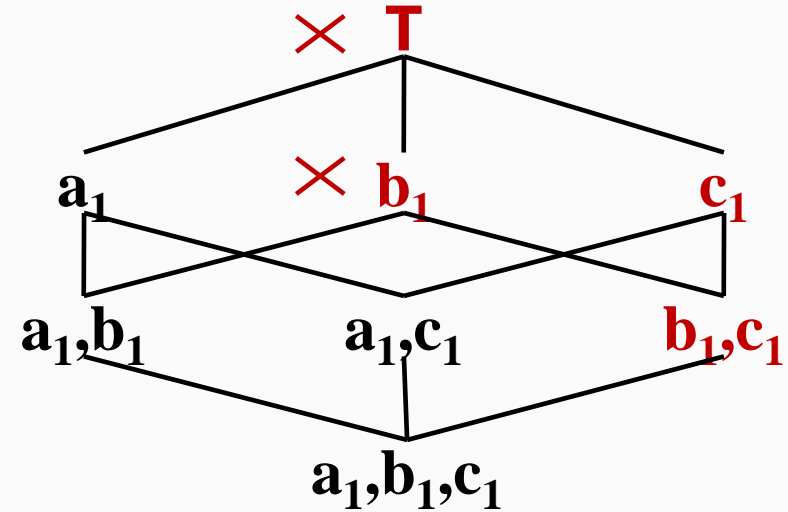
Brute-Force Approach

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



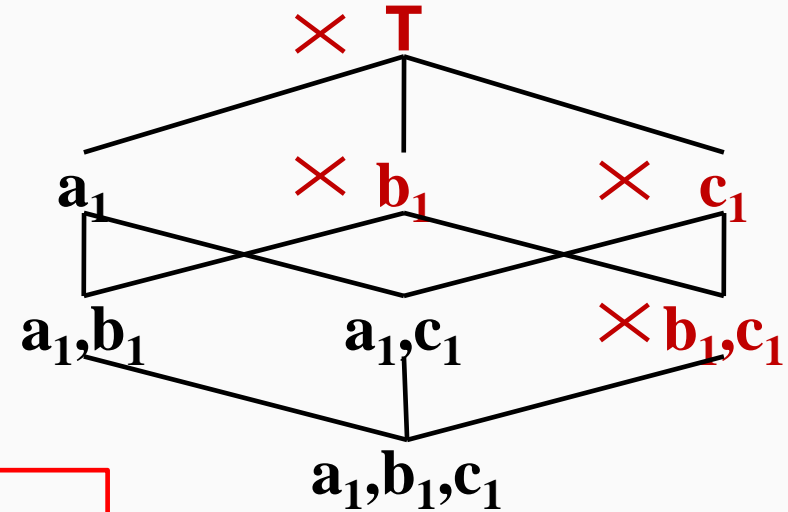
Brute-Force Approach

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



Brute-Force Approach

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



Total $|R| * (2^{|\mathcal{D}|+|\mathcal{M}|-1})$ comparisons!
 Total 16 comparisons in this case!

Challenges

- Exhaustive comparison with every tuple
- Under every constraint
- Over every measure subspace



Challenges and Ideas

➤ Exhaustive comparison with every tuple

✓ Tuple reduction

■ Comparison with skyline tuples is enough

■ $t_4 \succ_{\{m_1, m_2\}} t_3 \succ_{\{m_1, m_2\}} t_5 \Rightarrow t_4 \succ_{\{m_1, m_2\}} t_5$

<i>id</i>	<i>d₁</i>	<i>d₂</i>	<i>d₃</i>	<i>m₁</i>	<i>m₂</i>
<i>t₁</i>	<i>a₁</i>	<i>b₁</i>	<i>c₁</i>	10	5
<i>t₂</i>	<i>a₁</i>	<i>b₁</i>	<i>c₁</i>	15	10
<i>t₃</i>	<i>a₁</i>	<i>b₁</i>	<i>c₁</i>	17	17
<i>t₄</i>	<i>a₁</i>	<i>b₁</i>	<i>c₁</i>	20	20
<i>t₅</i>	<i>a₁</i>	<i>b₁</i>	<i>c₁</i>	11	15

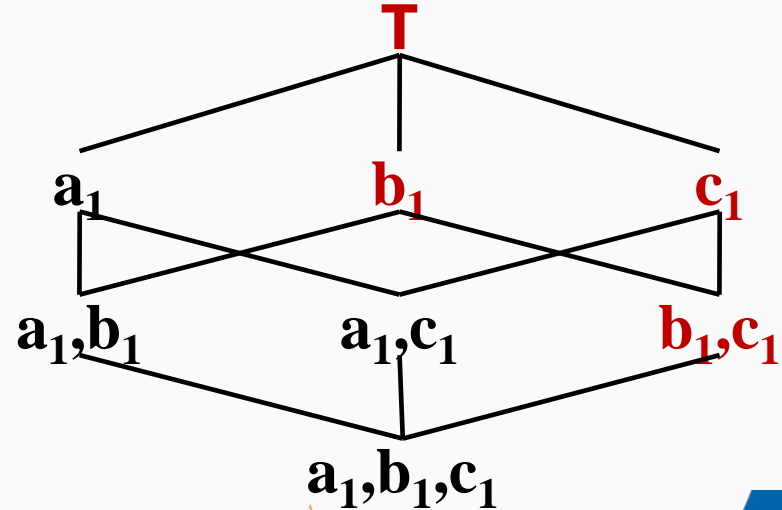
Challenges and Ideas

➤ Under every constraint

✓ Constraint pruning

■ In $C^{t,t'}$, one comparison on t and t' is enough

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



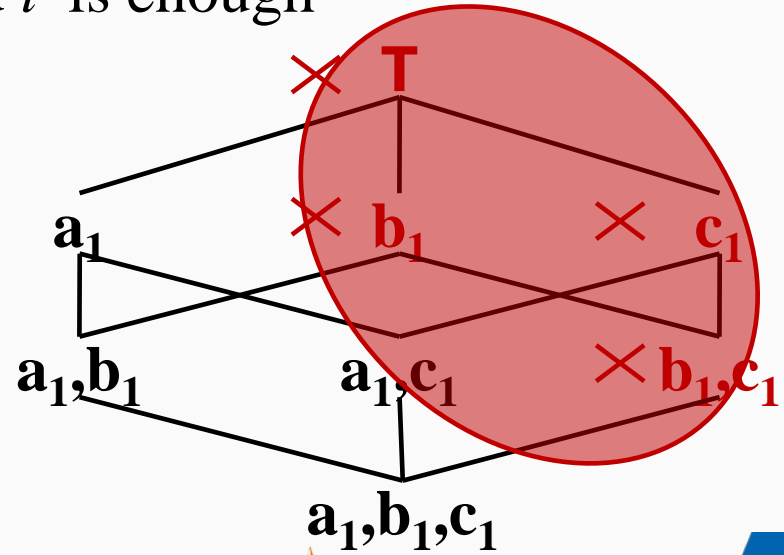
Challenges and Ideas

➤ Under every constraint

✓ Constraint pruning

■ In $C^{t,t'}$, one comparison on t and t' is enough

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



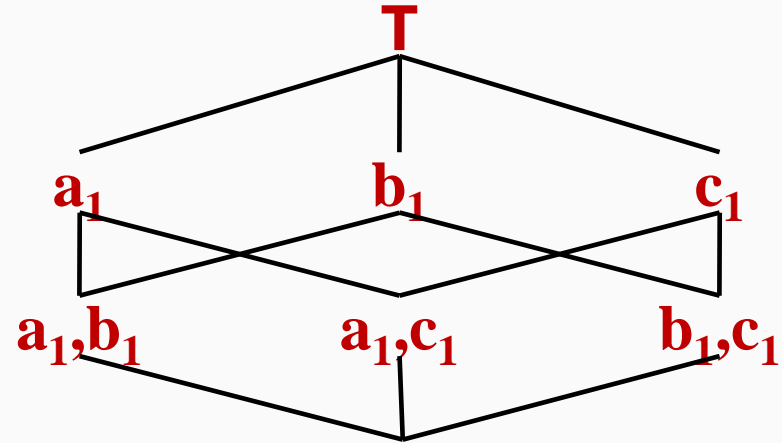
Challenges and Ideas

➤ Over every measure subspace

✓ Sharing computation across measure subspaces

■ Reusing computations on full space in subspaces

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



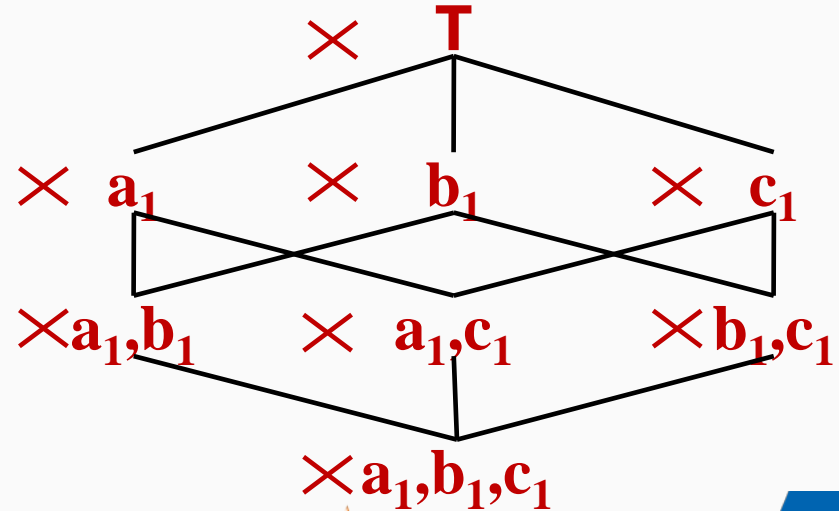
Challenges and Ideas

➤ Over every measure subspace

✓ Sharing computation across measure subspaces

■ Reusing computations on full space in subspaces

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	
t_2	a_1	b_1	c_1	15	
t_3	a_2	b_1	c_2	17	
t_4	a_2	b_1	c_1	20	
t_5	a_1	b_1	c_1	11	



Our Algorithms

- Tuple reduction + Constraint pruning
 - BottomUp
 - TopDown
- Tuple reduction + Constraint pruning + Sharing computation
 - SBottomUp
 - STopDown



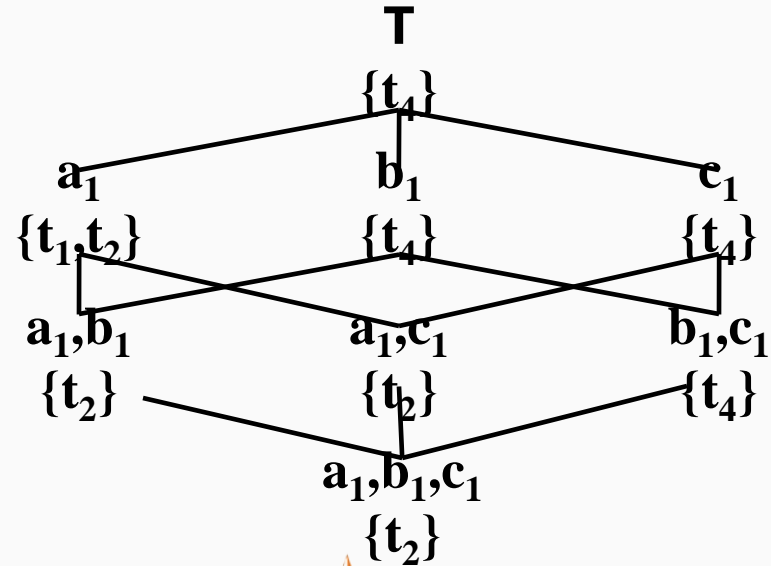
BottomUp

- Stores a tuple for every such constraint that qualifies it as a contextual skyline tuple
- Traverses the constraints in C^t in a bottom-up, breadth-first manner



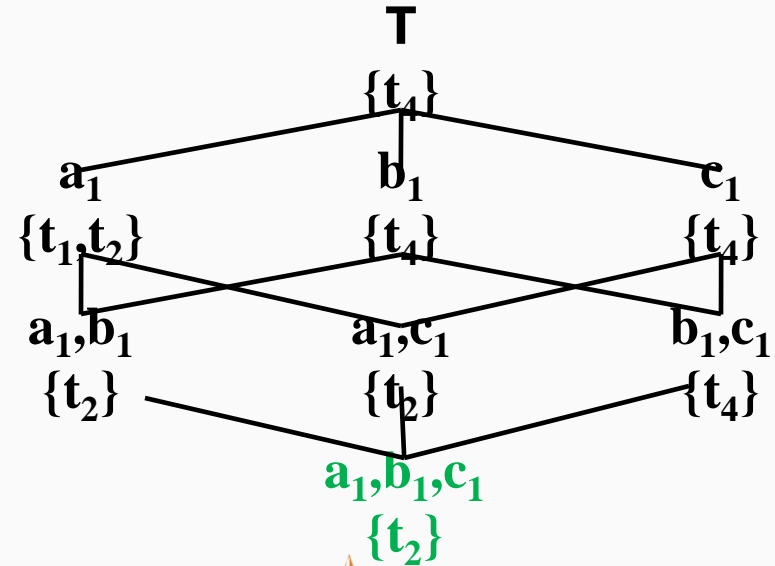
BottomUp

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



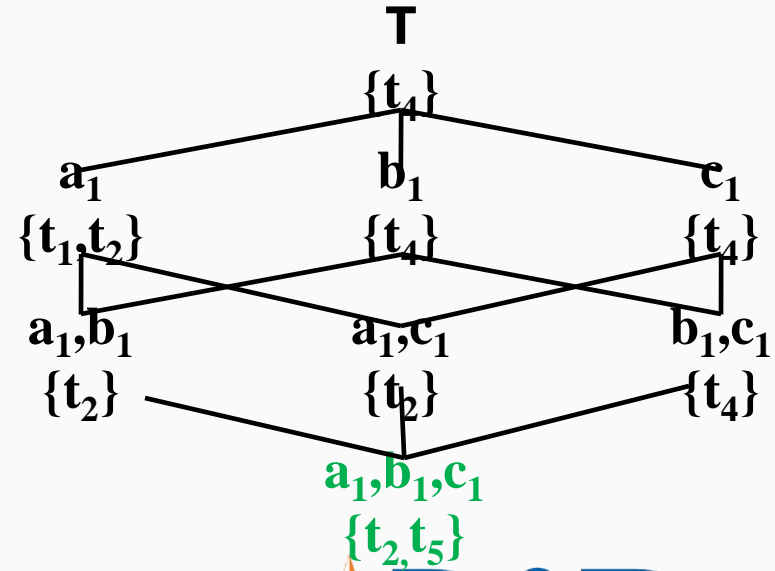
BottomUp

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_1	c_1	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_1	b_1	c_1	17	17
t_4	a_1	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



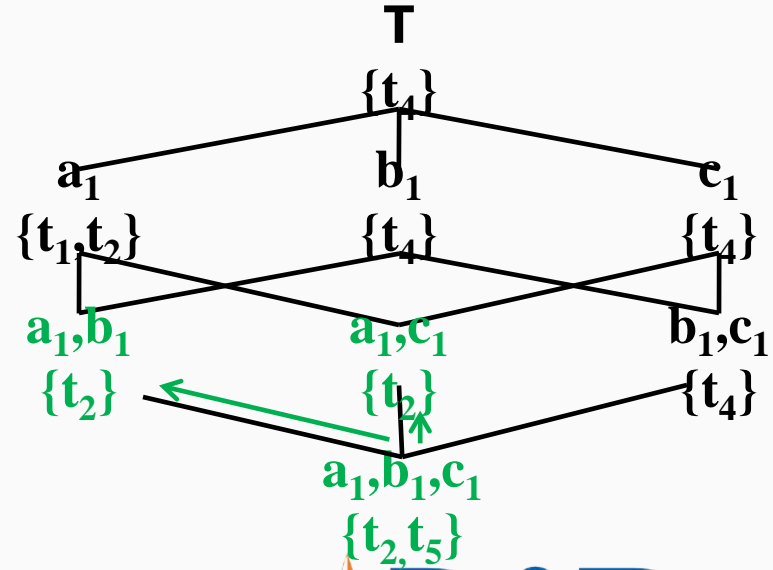
BottomUp

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_1	c_1	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_2	c_2	17	17
t_4	a_2	b_2	c_2	20	20
t_5	a_1	b_1	c_1	11	15



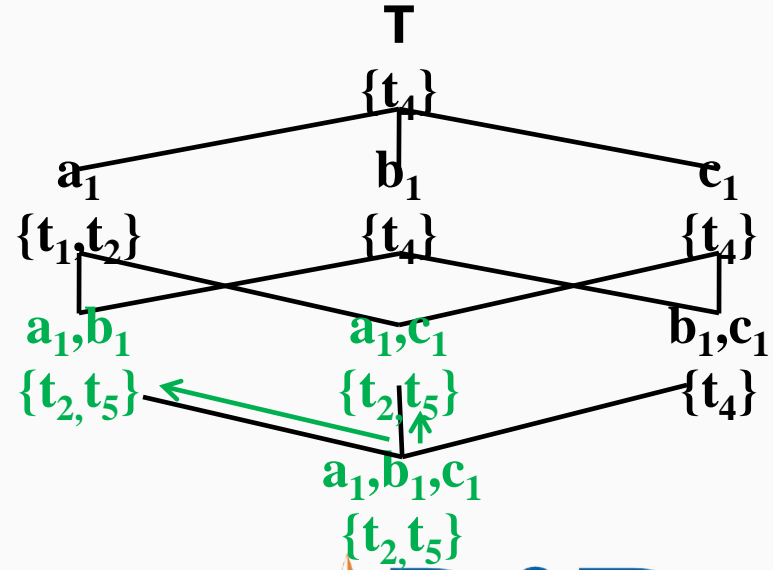
BottomUp

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_1	c_1	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_1	b_1	c_1	17	17
t_4	a_1	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



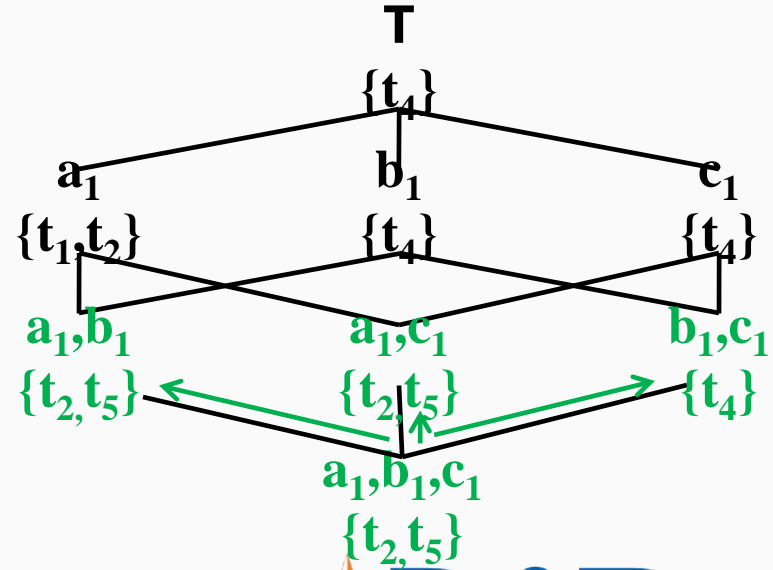
BottomUp

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_1	c_1	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_2	c_2	17	17
t_4	a_2	b_2	c_2	20	20
t_5	a_1	b_1	c_1	11	15



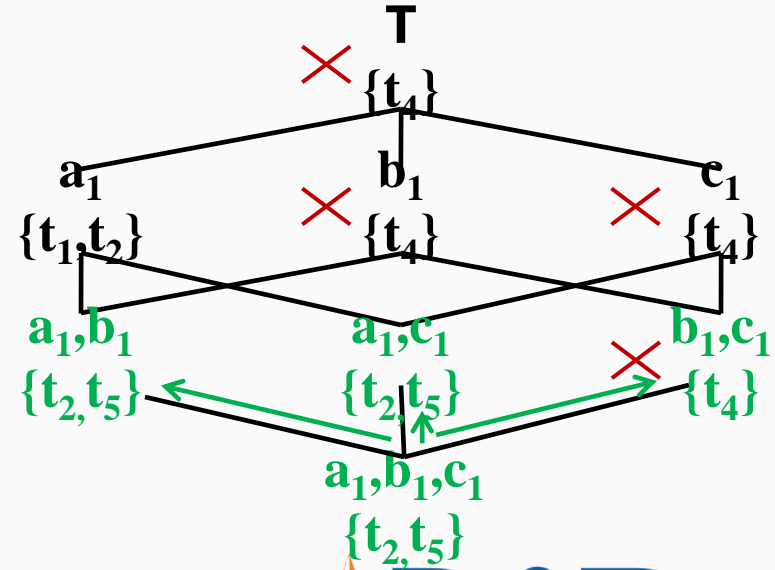
BottomUp

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	a_1	a_1	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_1	b_1	c_1	17	17
t_4	a_1	b_1	c_1	20	20
t_5		b_1	c_1	11	15



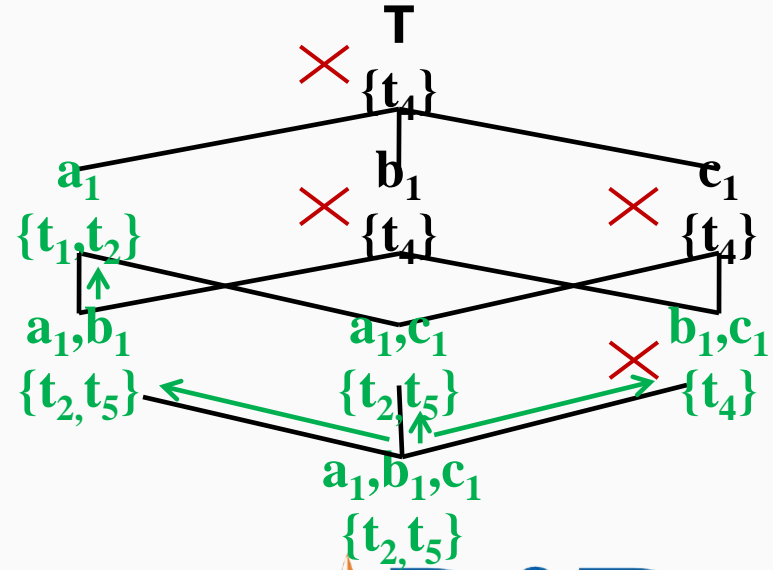
BottomUp

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	a_1	a_1	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_1	b_1	c_1	17	17
t_4	a_1	b_1	c_1	20	20
t_5		b_1	c_1	11	15



BottomUp

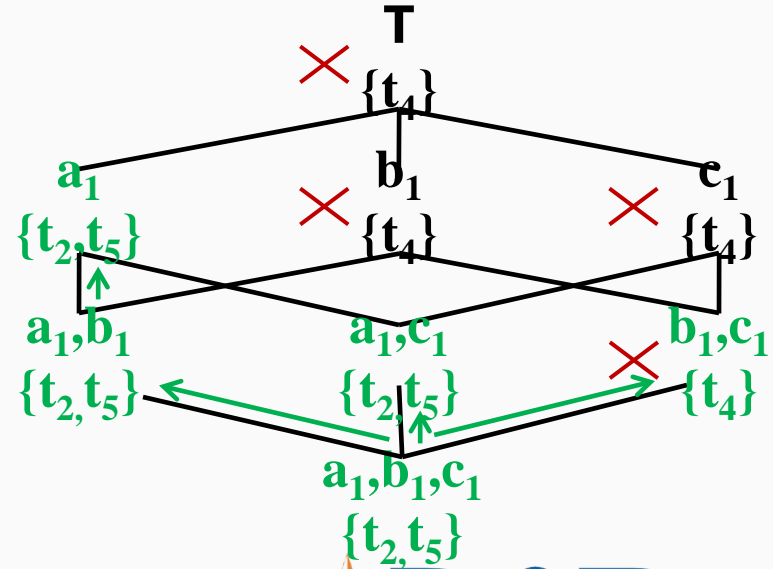
id	d_1	d_2	d_3	m_1	m_2
t_1	a_1			10	15
t_2	a_1			15	10
t_3		b_1		17	17
t_4		a_1	c_1	20	20
t_5	a_1			11	15



BottomUp

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_1	c_1	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_1	b_1	c_1	17	17
t_4	a_1	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15

6 comparisons in this case



BottomUp

➤ Cons of BottomUp

- Repetitive storage: space complexity
- Repetitive comparisons: time complexity

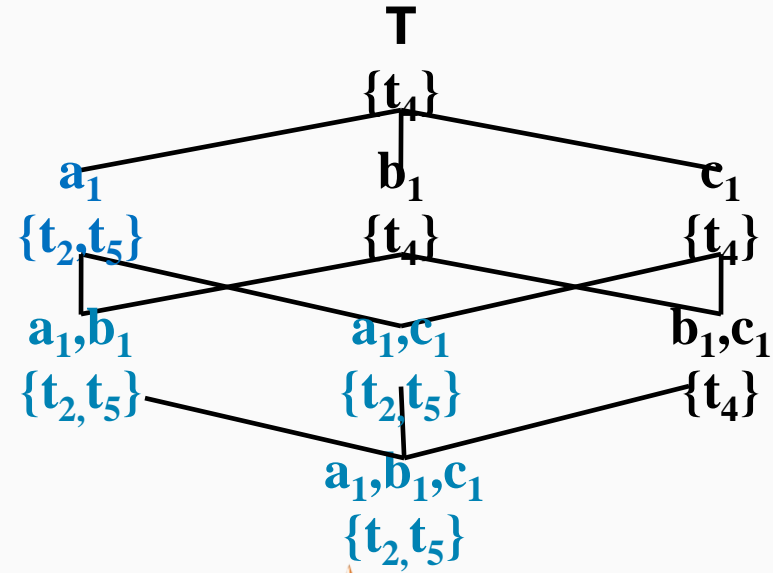
TopDown stores a tuple for its maximal skyline constraints only.



Skyline Constraints

Constraints whose contextual skylines include t .

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15

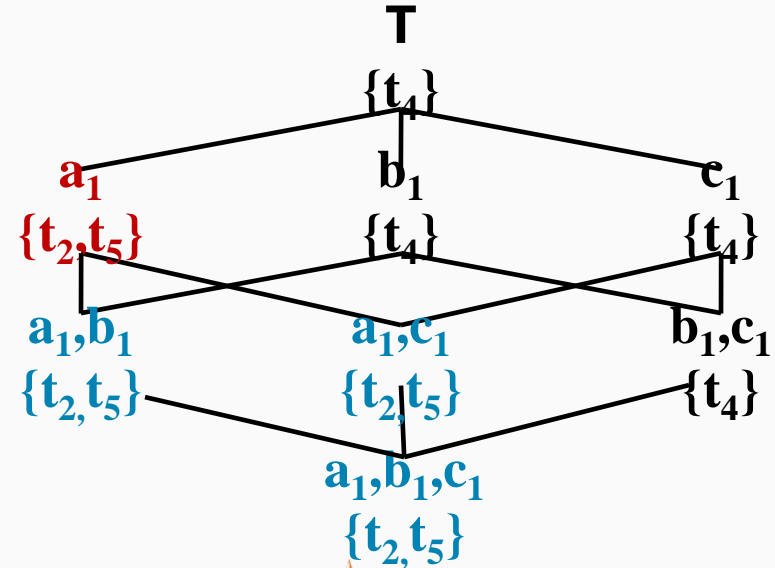


TopDown

Maximal Skyline Constraints

Constraints not subsumed by any other skyline constraints of t .

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15

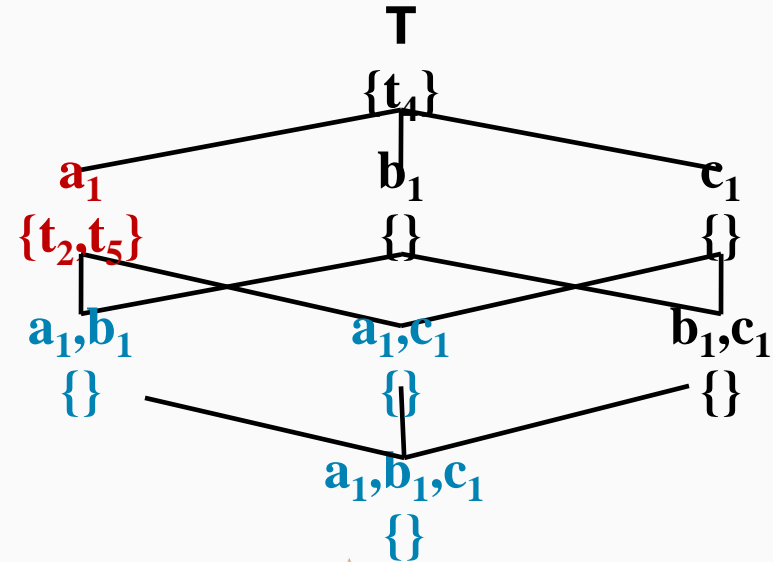


TopDown

Maximal Skyline Constraints

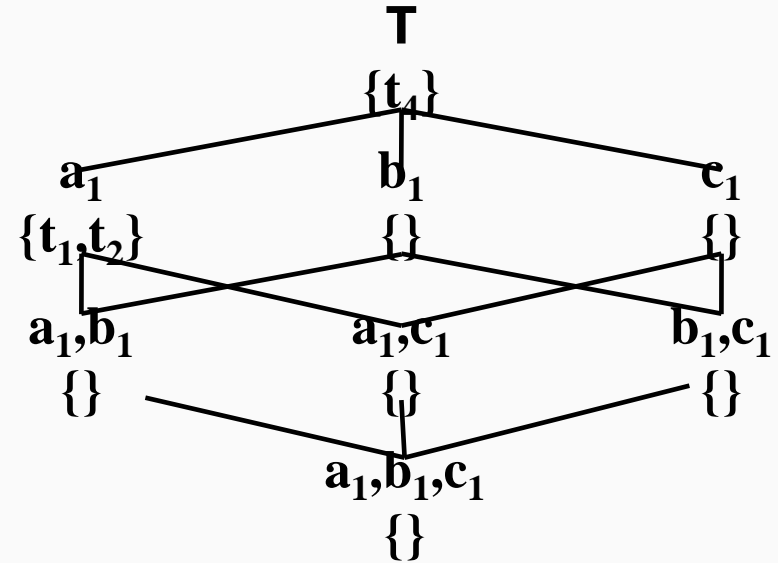
Constraints not subsumed by any other skyline constraints of t .

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



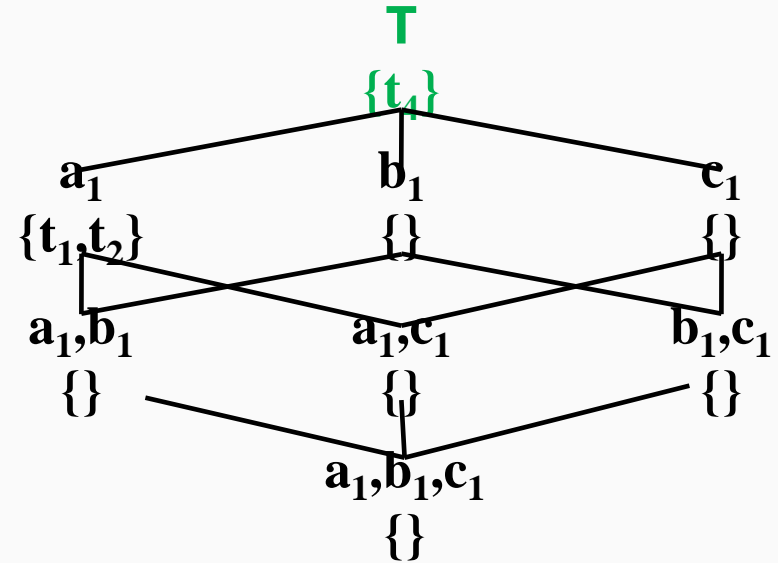
TopDown

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



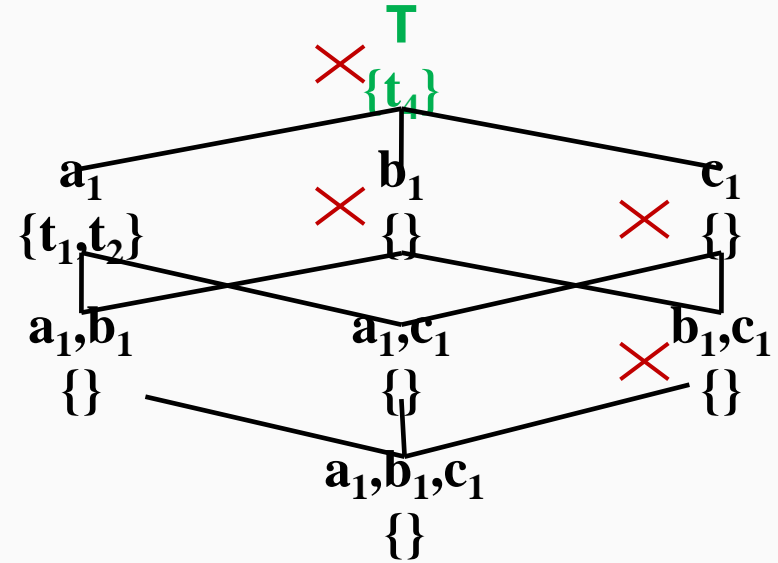
TopDown

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



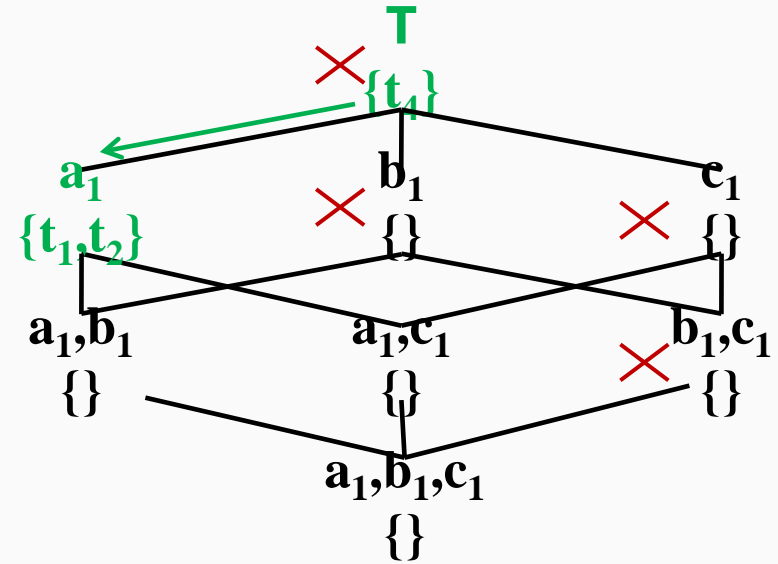
TopDown

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	17
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	15



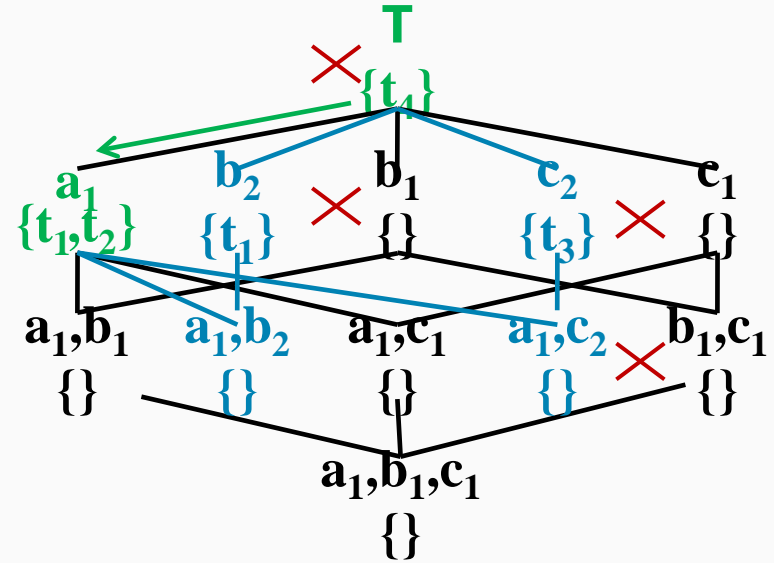
TopDown

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1			10	15
t_2	a_1			15	10
t_3		b_1		17	17
t_4		a_1	c_1	20	20
t_5	a_1			11	15



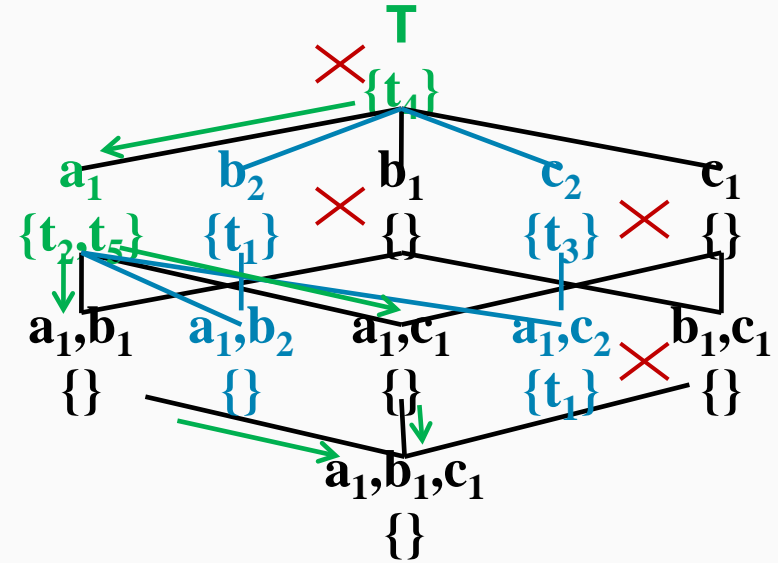
TopDown

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1			10	15
t_2	a_1			15	10
t_3		b_1		17	17
t_4		a_1	c_1	20	20
t_5	a_1			11	15



TopDown

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1			10	15
t_2	a_1			15	10
				17	17
				20	20
t_5	a_1			11	15



3 comparisons in this case

STopDown and SBottomUp

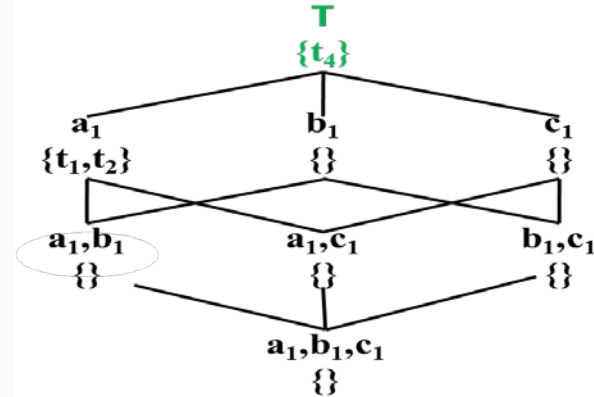
➤ Con of BottomUp and TopDown

- Need to compute **over every measure subspace** separately
 - STopDown and SBottomUp share computation across different subspaces



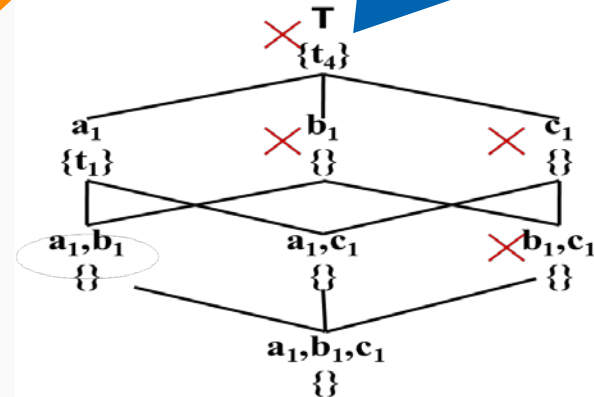
STopDown

<i>id</i>	<i>d</i> ₁	<i>d</i> ₂	<i>d</i> ₃	<i>m</i> ₁	<i>m</i> ₂
<i>t</i> ₁	<i>a</i> ₁	<i>b</i> ₂	<i>c</i> ₂	10	15
<i>t</i> ₂	<i>a</i> ₁	<i>b</i> ₁	<i>c</i> ₁	15	10
<i>t</i> ₃	<i>a</i> ₂	<i>b</i> ₁	<i>c</i> ₂	17	17
<i>t</i> ₄	<i>a</i> ₂	<i>b</i> ₁	<i>c</i> ₁	20	20
<i>t</i> ₅	<i>a</i> ₁	<i>b</i> ₁	<i>c</i> ₁	11	15



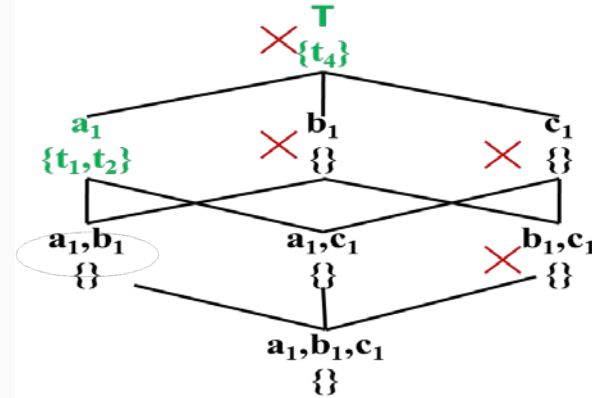
Comparison with *t*₄ is skipped

<i>id</i>	<i>d</i> ₁	<i>d</i> ₂	<i>d</i> ₃	<i>m</i> ₁	<i>m</i> ₂
<i>t</i> ₁	<i>a</i> ₁	<i>b</i> ₂	<i>c</i> ₂	10	15
<i>t</i> ₂	<i>a</i> ₁	<i>b</i> ₁	<i>c</i> ₁	15	10
<i>t</i> ₃	<i>a</i> ₂	<i>b</i> ₁	<i>c</i> ₂	17	17
<i>t</i> ₄	<i>a</i> ₂	<i>b</i> ₁	<i>c</i> ₁	20	20
<i>t</i> ₅	<i>a</i> ₁	<i>b</i> ₁	<i>c</i> ₁		15



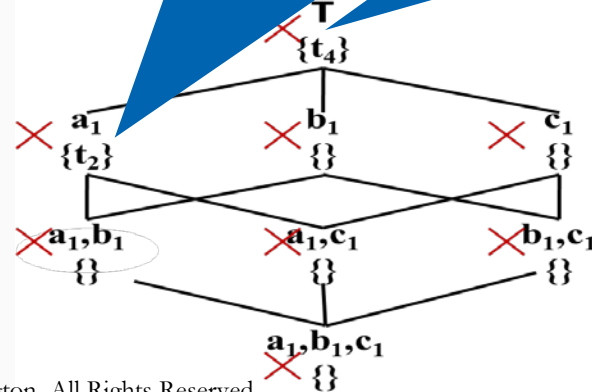
STopDown

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1			10	15
t_2	a_1			15	10
				20	20
t_5	a_1			11	15



Comparisons with t_2 & t_4 are skipped

id	d_1	d_2	d_3	m_1	m_2
t_1	a_1	b_2	c_2	10	15
t_2	a_1	b_1	c_1	15	10
t_3	a_2	b_1	c_2	17	
t_4	a_2	b_1	c_1	20	20
t_5	a_1	b_1	c_1	11	



Experiment Setup

❑ NBA Dataset

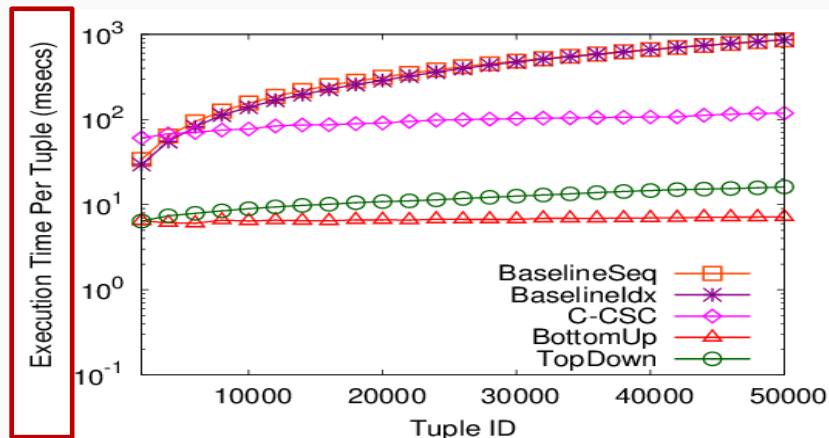
- 317,371 tuples of NBA box scores from 1991-2004 seasons
- 8 dimension attributes
- 7 measure attributes

❑ Weather Dataset

- 7.8 million tuples of weather forecast from different locations of six countries & regions of UK
- 7 dimension attributes
- 7 measure attributes



Memory-Based Implementation



NBA Dataset

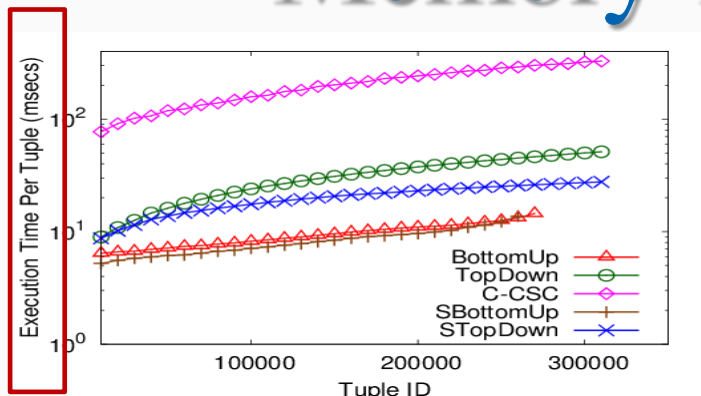
❑ Maintaining CSC for each constraint causes overhead

(Xia et al. SIGMOD 2006)

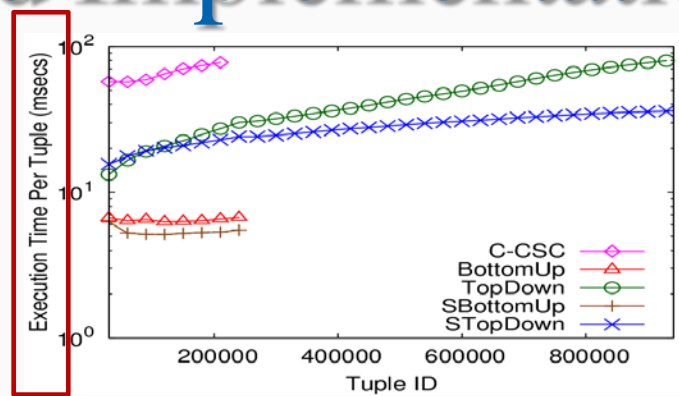
- Can't take advantage of constraint pruning



Memory-Based Implementation



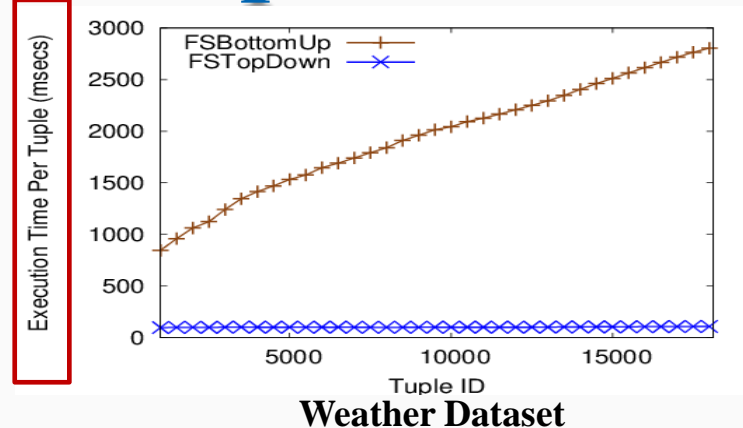
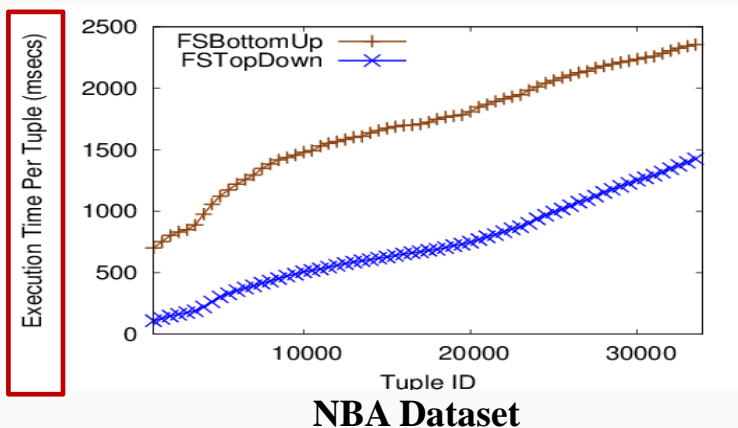
NBA Dataset



Weather Dataset

- ❑ BottomUp/SBottomUp exhausted available JVM heap
 - memory overflow
- ❑ TopDown / STopDown was outperformed by BottomUp/SBottomUp
 - Updating maximal skyline constraints causes overhead

File-Based Implementation



- ❑ Each (C, M) is stored in a binary file
- ❑ While traversing, file-read operation occurs if file is non-empty: **FSTopDown encounters many empty files**
- ❑ For updating, file-write operation occurs: **FSTopDown stores fewer tuples**
- ❑ I/O-cost dominates in-memory computation

Discovered Facts

- Lamar Odom had 30 points, 19 rebounds and 11 assists on March 6, 2004. No one before had a better or equal performance in NBA history.
- Allen Iverson had 38 points and 16 assists on April 14, 2004 to become the first player with a 38/16 (points/assists) game in the 2004-2005 season.
- Damon Stoudamire scored 54 points on January 14, 2005. It is the highest score in history made by any Trail Blazers.



Prominent Streak Discovery in Sequence Data. [Xiao Jiang, Chengkai Li, Ping Luo, Min Wang, Yong Yu. KDD 2011, pages 1280-1288.](#)

Discovering General Prominent Streaks in Sequence Data.
[Gensheng Zhang, Xiao Jiang, Ping Luo, Min Wang, Chengkai Li. ACM TKDD, 8\(2\):article 9, June 2014.](#)



Prominent Streaks

Prominent streaks stated in news articles:

"This month the Chinese capital has experienced 10 days with a maximum temperature in around 35 degrees Celsius – the most for the month of July in a decade."

"The Nikkei 225 closed below 10000 for the 12th consecutive week, the longest such streak since June 2009."

"He (LeBron James) scored 35 or more points in nine consecutive games and joined Michael Jordan and Kobe Bryant as the only players since 1970 to accomplish the feat."



Concepts

Streak

Input: a sequence of values

Streak $\langle [l, r], v \rangle$ is a triple: left-end (l), right-end (r), minimum value in interval $[l, r]$

3 1 7 7 2 5 4 6 7 3
 $\langle [6, 8], 4 \rangle$

Streak dominance relation

$s1 = \langle [l1, r1], v1 \rangle$ dominates $s2 = \langle [l2, r2], v2 \rangle$ iff
 $r1 - l1 > r2 - l2, v1 \geq v2$ or $r1 - l1 \geq r2 - l2, v1 > v2$

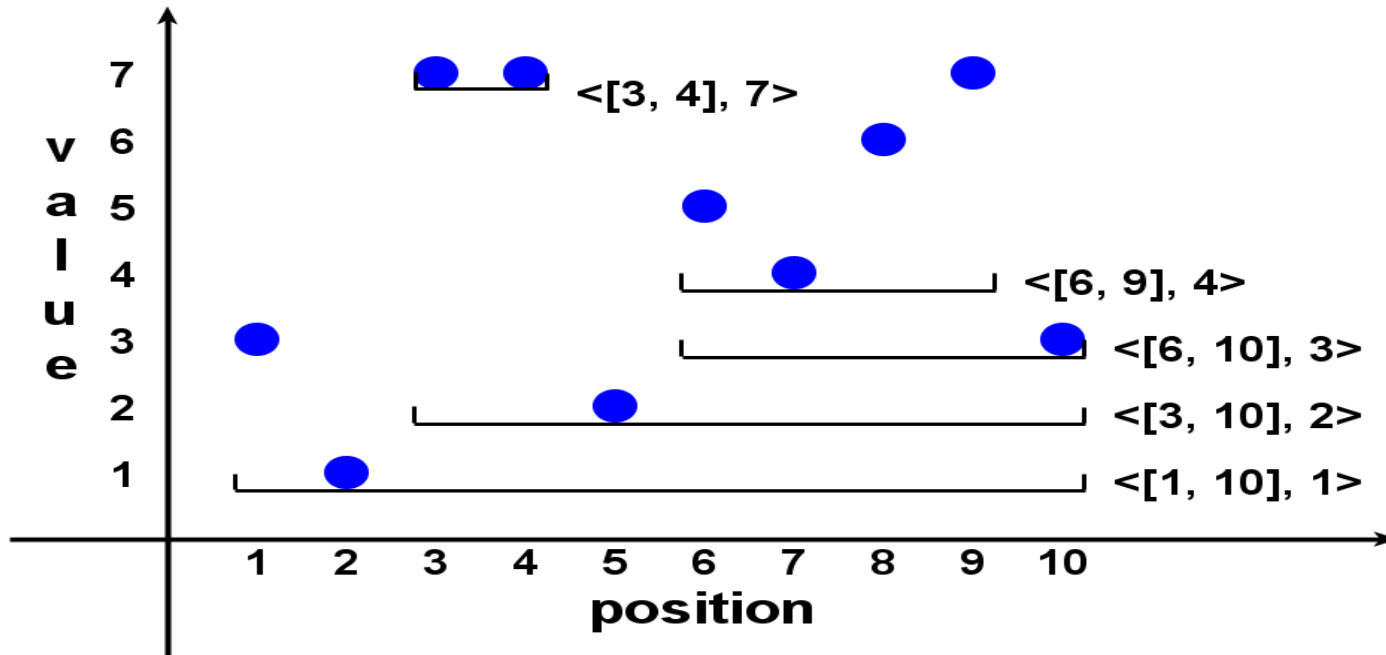
Prominent streaks (PS)

A streak is prominent if it is not dominated by any other streaks.



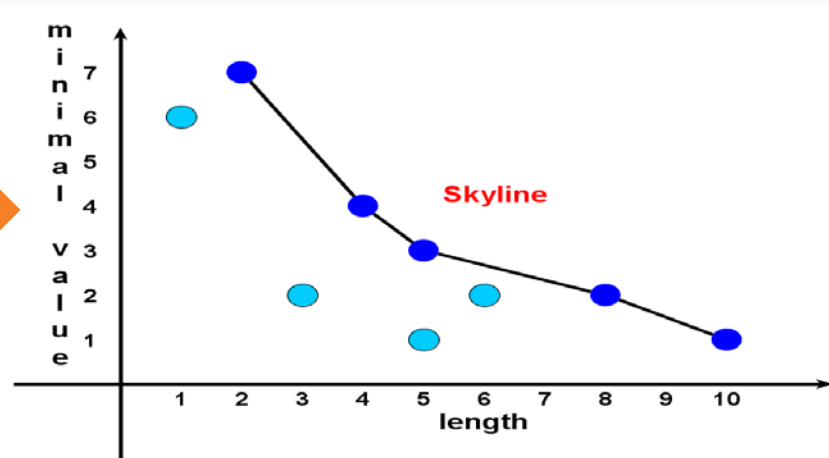
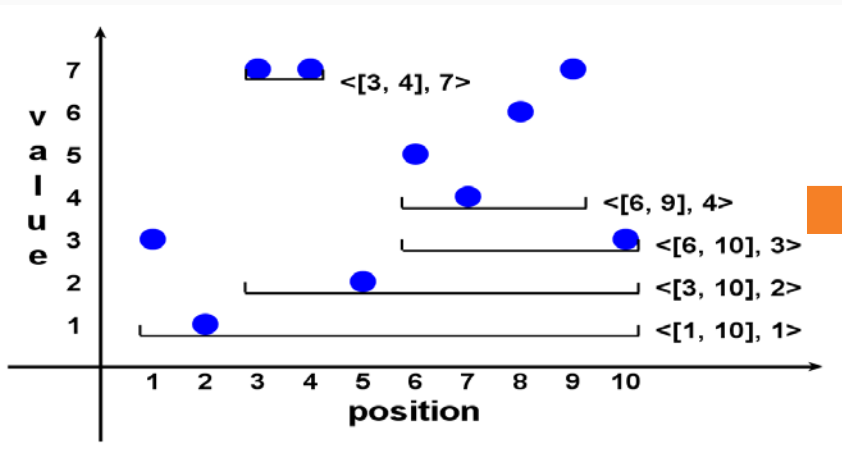
Example

3 1 7 7 2 5 4 6 7 3



Prominent Streaks are Skyline Points in 2-d Space

3 1 7 7 2 5 4 6 7 3



Tasks

Task 1: discovery

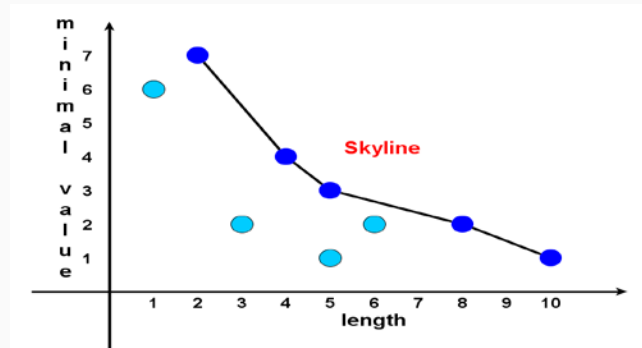
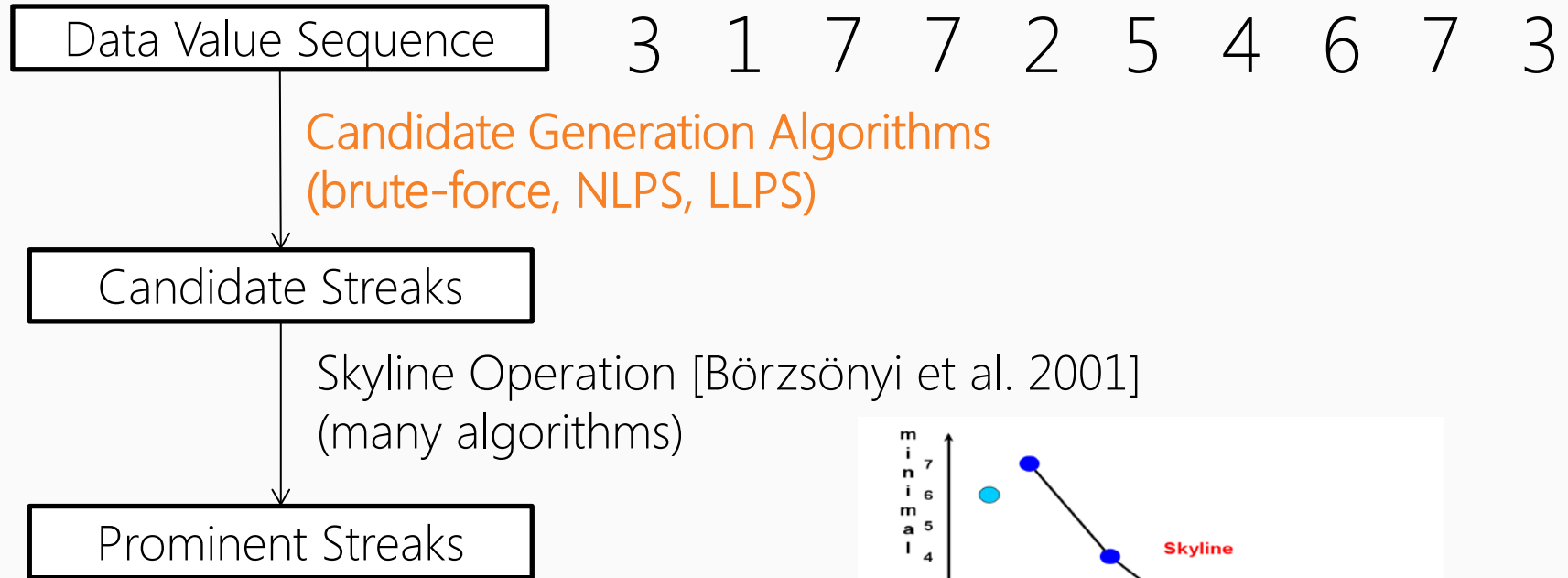
Find all prominent streaks in a sequence

Task 2: monitoring

Always keep prominent streaks up-to-date, when sequence grows (real-world sequences often grow)



Solution Framework



Candidate Generation: Number Of Candidates

Brute-force

Quadratic

NLPS

Superlinear

LLPS

Linear



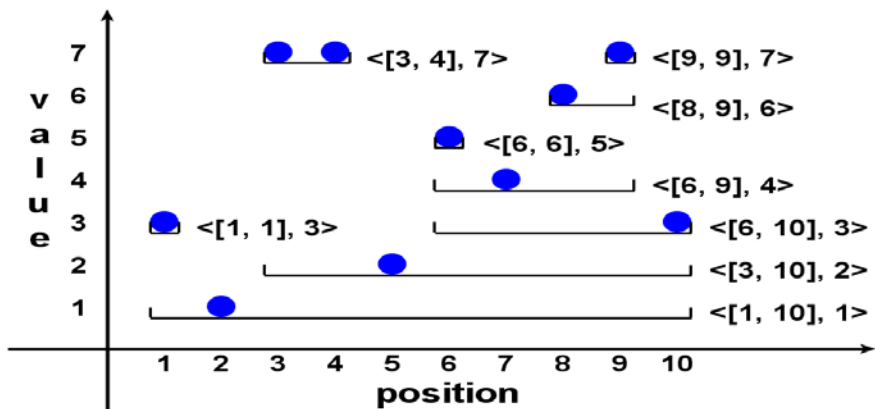
Local Prominent Streak

Local dominance relation

$s1 = \langle [l1, r1], v1 \rangle$ locally dominates $s2 = \langle [l2, r2], v2 \rangle$ iff
 $s1$ dominates $s2$ and $[l1, r1] \supset [l2, r2]$

Local prominent streak (LPS)

A streak is locally prominent if it is not locally dominated by any other streaks.



Important Properties

(1) LPS is sufficient

A prominent streak must be an LPS.

(2) LPS is small

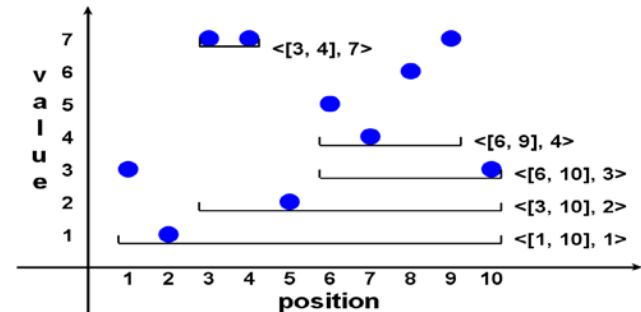
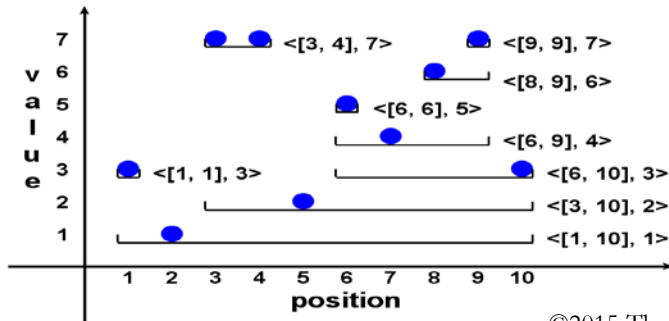
The number of LPSs is less than or equal to the sequence length.

(Hint: The number of LPSs getting min value at position k is at most 1.)

Conclusion

LPS is an excellent set of candidate streaks, of linear size.

Candidate generation problem \Rightarrow finding local prominent streaks



Linear LPS (LLPS) Method

Sequence p_1, p_2, \dots, p_n .

1. Maintain a list of candidate streaks when scanning the sequence rightward.
2. After p_k , right-ends of candidates are all k .
3. At p_{k+1} , try to extend the candidates rightward.

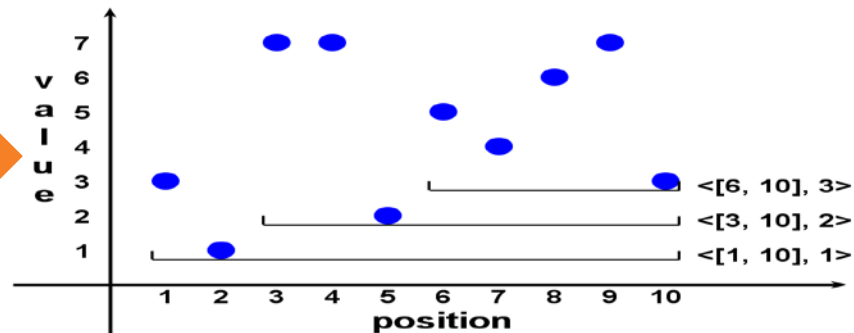
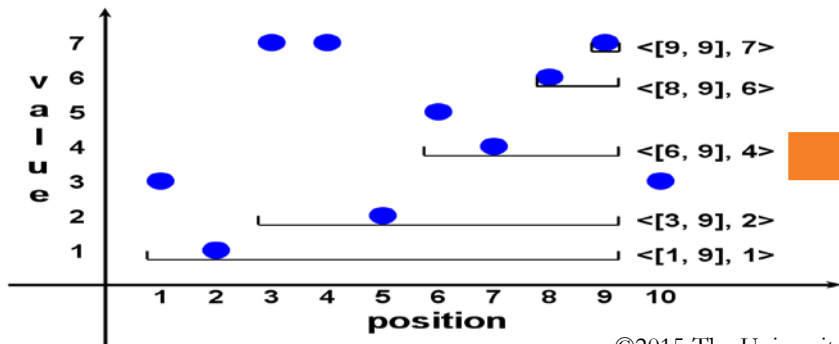
Candidates s :

(3.a) $s.v < p_{k+1}$: extend.

(3.b) $s.v > p_{k+1}$: belong to LPS.

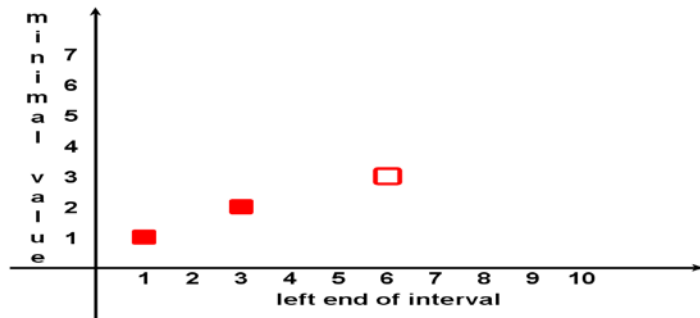
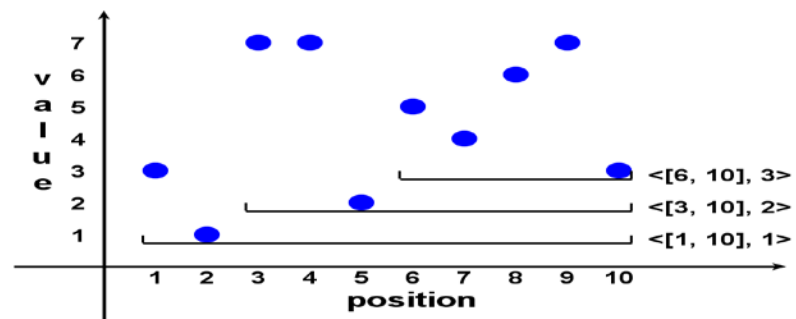
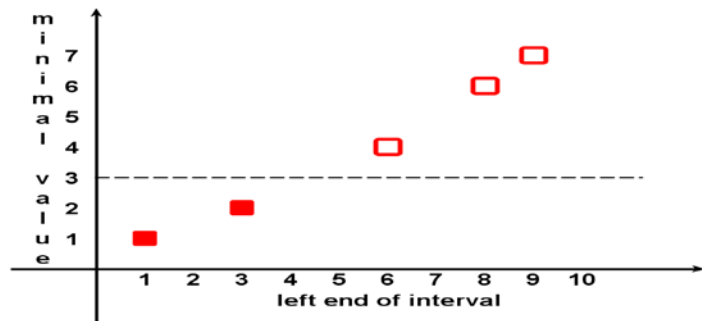
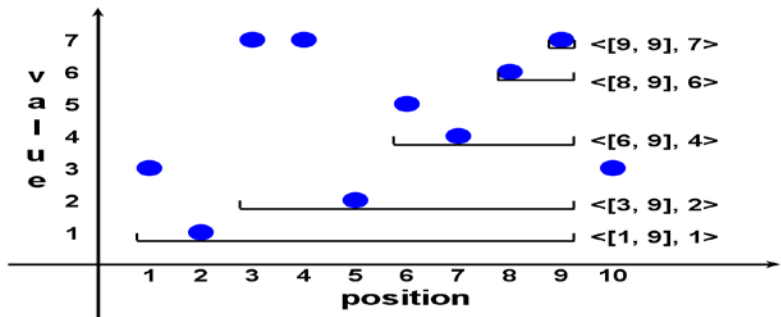
(3.c) $s.v \geq p_{k+1}$: extend the leftmost (longest) such s .

4. After p_n all remaining candidates are LPS.



Linear LPS (LLPS) Method

Candidates share the same right-end, their minimum values monotonically increase, if they are listed in the increasing order of left-ends.



Linear LPS (LLPS) Method

After p_k , it has found:

All LPSs ending before k

Candidates ending at k either are LPSs or can be grown to LPSs ending after k .

Monitoring (keeping prominent streaks up-to-date) is simple:

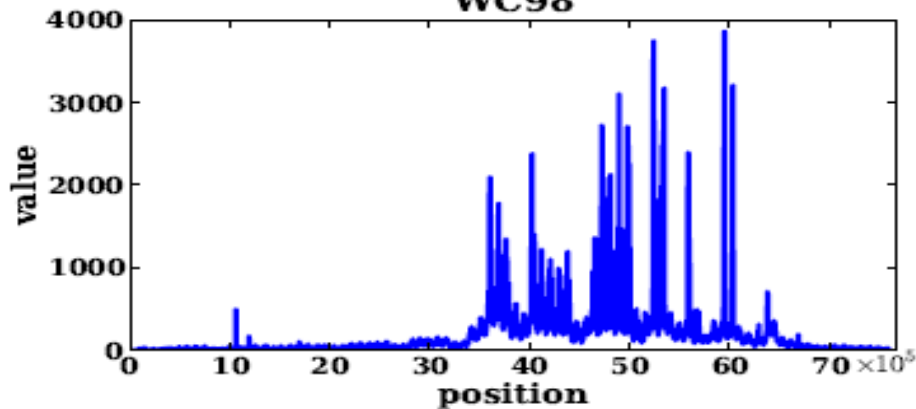
If PSs till k are requested, compare all found LPSs and all remaining candidates.



Datasets In Experiments

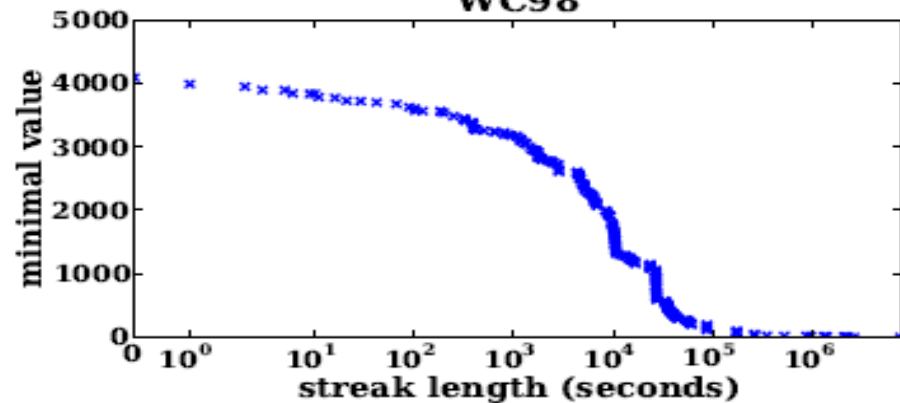
name	length	# prominent streaks	description
Gold	1074	137	Daily morning gold price in US dollars, 01/1985-03/1989.
River	1400	93	Mean daily flow of Saugeen River near Port Elgin, 01/1988-12/1991.
Melb1	3650	55	The daily minimum temperature of Melbourne, Australia, 1981-1990.
Melb2	3650	58	The daily maximum temperature of Melbourne, Australia, 1981-1990.
Wiki1	4896	58	Hourly traffic to en.wikipedia.org/wiki/Main_page , 04/2010-10/2010.
Wiki2	4896	51	Hourly traffic to en.wikipedia.org/wiki/Lady_gaga , 04/2010-10/2010.
Wiki3	4896	118	Hourly traffic to en.wikipedia.org/wiki/Inception_(film) , 04/2010-10/2010.
SP500	10136	497	S&P 500 index, 06/1960-06/2000.
HPQ	12109	232	Closing price of HPQ in NYSE for every trading day, 01/1962-02/2010.
IBM	12109	198	Closing price of IBM in NYSE for every trading day, 01/1962-02/2010.
AOL	132480	127	Number of queries sent to AOL search engine in every minute over three months.
WC98	7603201	286	Number of requests to World Cup 98 web site in every second, 04/1998-07/1998.

WC98



(a) Data Sequence

WC98



(b) Prominent Streaks

Sample Prominent Streaks

Melbourne daily min/max temperature between 1981 and 1990 (Melb1 & Melb2)

More than 2000 days with min temperature above zero
6 days: the longest streak above 35 degrees Celsius



Traffic count of Wikipedia page of Lady Gaga (Wiki2)

More than half of the prominent streaks are around Sep. 12th (VMA 2010)
at least 2000 hourly visits lasting for almost 4 days



General Prominent Streaks

Top-k, multi-dimensional and multi-sequence PS

"He (LeBron James) scored 35 or more points in nine consecutive games and joined Michael Jordan and Kobe Bryant as the only players since 1970 to accomplish the feat."

"Only player in NBA history to average at least 20 points, 10 rebounds and 5 assists per game for 6 consecutive seasons." (http://en.wikipedia.org/wiki/Kevin_Garnett)

NLPS/LLPS extended to such general PSs



Experiments On Multi-Sequence PSs

Table IX. Multi-sequence Prominent Streaks in Datas NBA1.

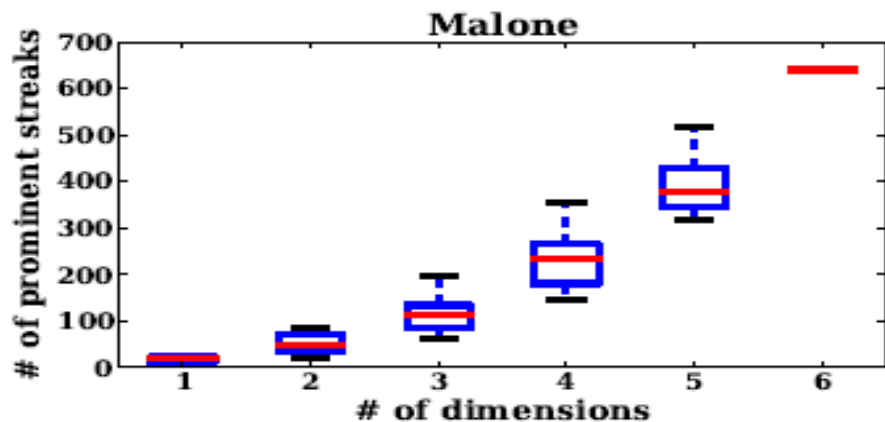
length	minimal value	players
1	71	David Robinson
2	51	Allen Iverson; Antawn Jamison
4	42	Kobe Bryant
9	40	Kobe Bryant
13	35	Kobe Bryant
14	32	Kobe Bryant
16	30	Kobe Bryant
17	27	Michael Jordan
27	26	Allen Iverson
34	24	Tracy McGrady
45	21	Allen Iverson
57	20	Allen Iverson
74	19	Shaquille O'Neal
94	18	Shaquille O'Neal
96	17	Karl Malone
119	16	Karl Malone
149	15	Karl Malone
159	14	Karl Malone
263	13	Karl Malone
357	12	Karl Malone
527	11	Karl Malone
575	10	Karl Malone
758	7	Karl Malone
858	6	Shaquille O'Neal
866	2	Karl Malone
932	1	John Stockton
1185	0	Jim Jackson



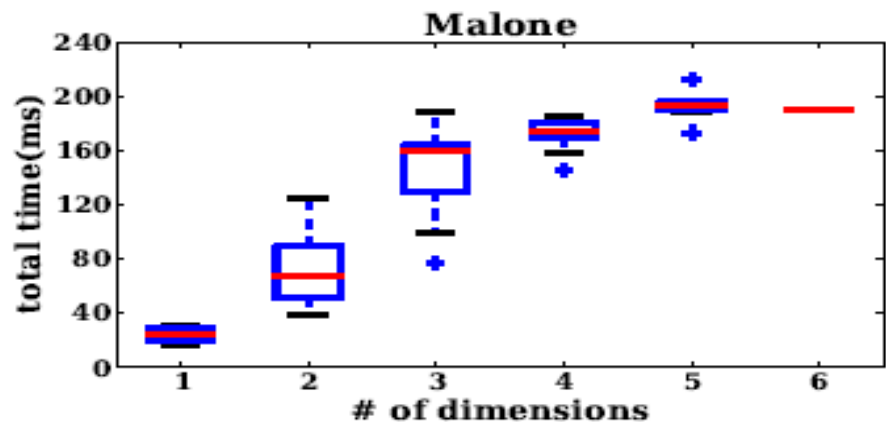
Experiments On Multi-Dim PSs

Table X. Data Sequences Used in Experiments on Multi-dimensional Prominent Streak Discovery.

name	length	# prominent streaks	# dimensions	description
Malone	986	640	6	1991-2004 game log of Karl Malone (minutes, points, rebounds, assists, steals, blocks)



(a) Number of Prominent Streaks



(b) Execution Time of LLPS

Fig. 13. Experiments on Increasing Dimensionality.

Experiments On General PSs

Table XIII. Data Sequences Used in Experiments on Top-5 Multi-sequence Multi-dimensional Prominent Streak Discovery.

name	# sequences	average length	# dimensions	# prominent streaks	description
NBA2	1185	290	6	10867	1991-2004 game log of all NBA players (minutes, points, rebounds, assists, steals, blocks)

Table XIV. Number of Candidate Streaks, Top-5 Multi-sequence Multi-dimensional Prominent Streak Discovery.

name	Baseline	NLPS	LLPS
NBA2	9.41×10^7	2.98×10^6	8.76×10^5

Table XV. Execution Time (in Milliseconds), Top-5 Multi-sequence Multi-dimensional Prominent Streak Discovery.

name	Baseline	NLPS	LLPS
NBA2	1.39×10^7	4.33×10^5	1.14×10^5

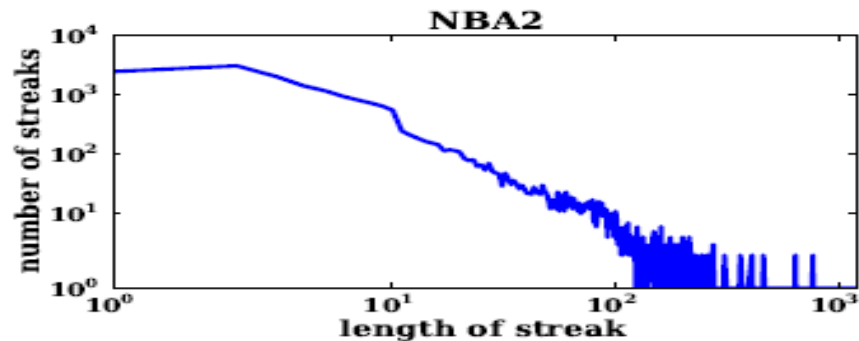


Fig. 14. Distribution of Prominent Streaks by Length.



On "One of the Few" Objects. You Wu, Pankaj K. Agarwal,
Chengkai Li, Jun Yang, Cong Yu. KDD 2012, pages 1487-1495



One-Of-The-Few Claims

Do these claims really hold water?

Karl Malone is **ONE OF THE ONLY TWO** players in NBA history with 25,000 points, 12,000 rebounds, and 5,000 assists in one's career.

He is **ONE OF THE ONLY THREE** candidates who have raised more than 25% from PAC contributions and 25% from self-financing.

How do we find truly interesting claims or individuals?



$X \text{ Is One-Of-}K \rightarrow X \text{ Is In } K\text{-Skyband}$

Claim

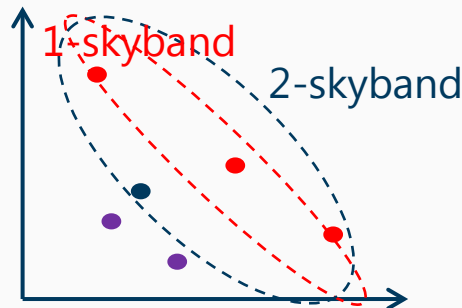
Karl Malone is **ONE OF THE ONLY TWO** players in NBA history with 25,000 points, 12,000 rebounds, and 5,000 assists in one's career.

General claim

Fewer than k objects dominate X in subspace of attributes $S \subseteq \{A_1, A_2, \dots, A_d\}$

k -skyband [Papadias et al. 2005] in S is the set of points each dominated by fewer than k other points in S

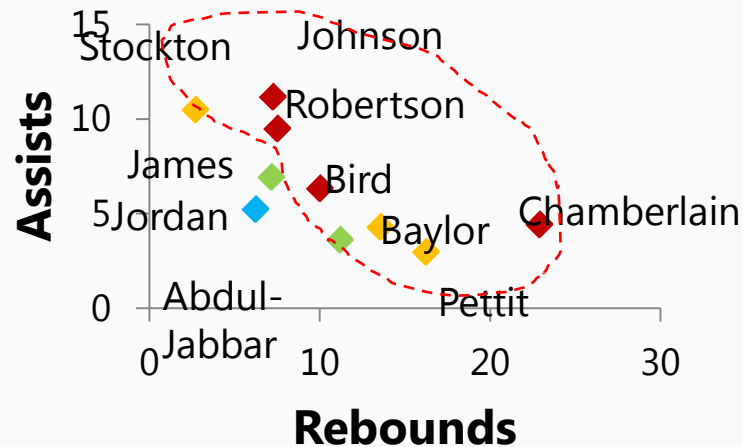
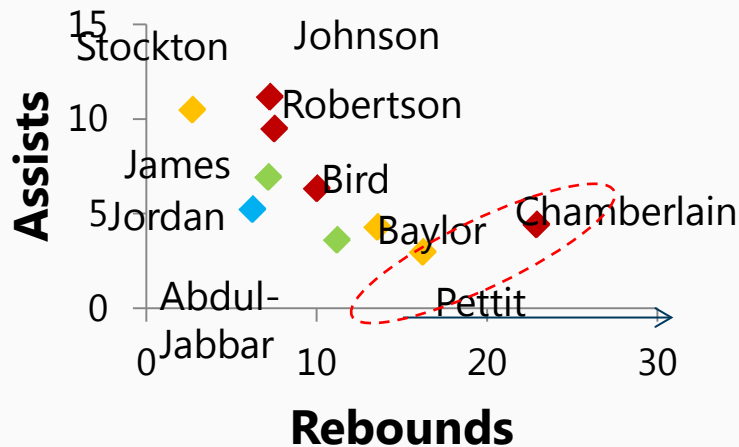
1-skyband : skyline



Small $K \neq$ Interesting

Subspaces are different

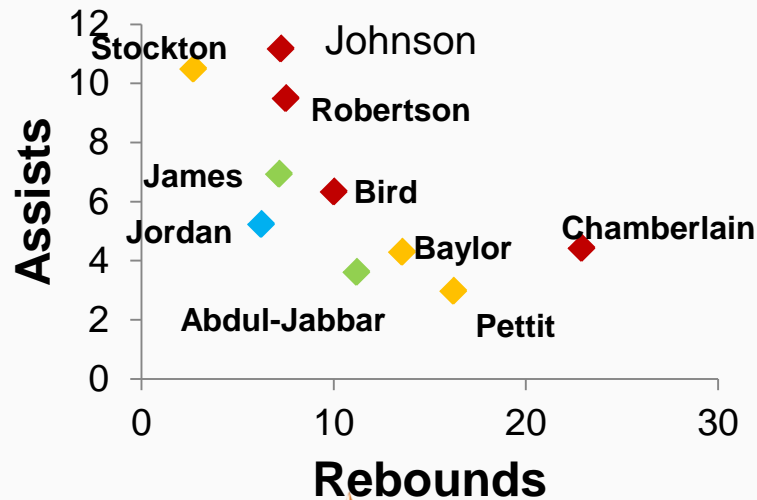
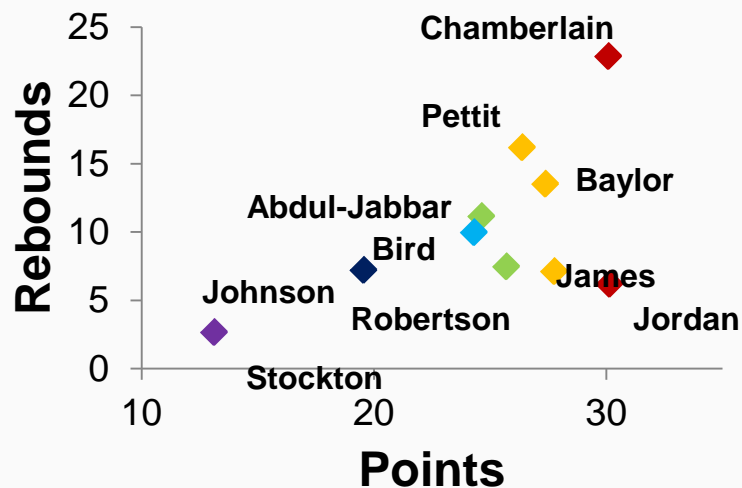
E.g., 2-sky and in {rebounds} vs. in {rebounds, assists}



Small $K \neq$ Interesting

Data distribution matters

E.g., 2-sky and in {points, rebounds} vs. in {rebounds, assists}



Top- τ Skyband

k -Skyband

Using the same k for all subspaces doesn't work

Asking user pick k for each subspace is infeasible

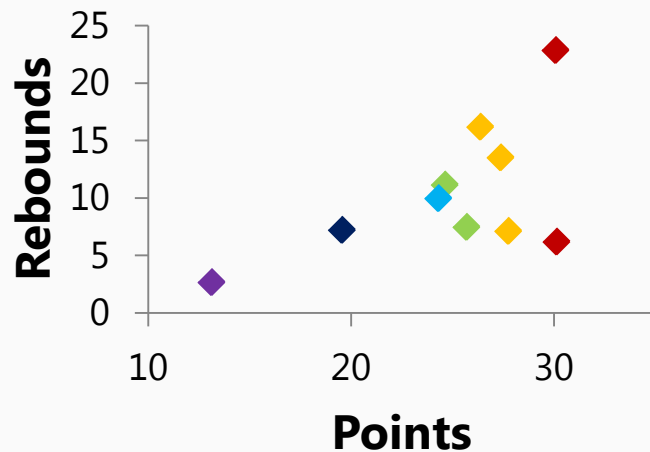
Top- τ Skyband

- User specifies a single parameter τ to cap # skyband objects.
- For each subspace \mathcal{S} , find its top- τ skyband, i.e., the largest k -skyband containing no more than τ objects

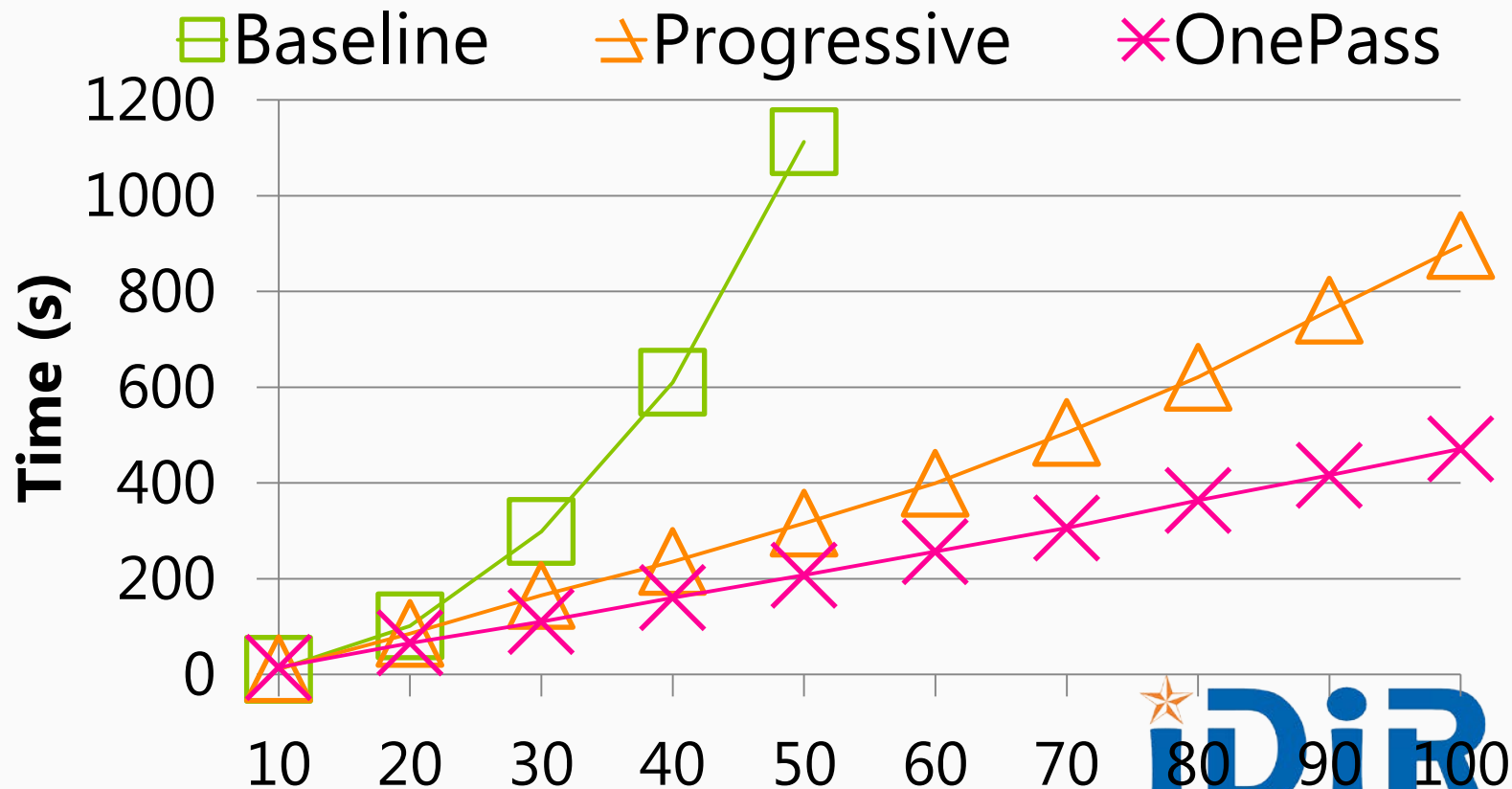
○ E.g., in {points, rebounds}:

$\tau=2 \rightarrow$ 1-skyband (size 2)

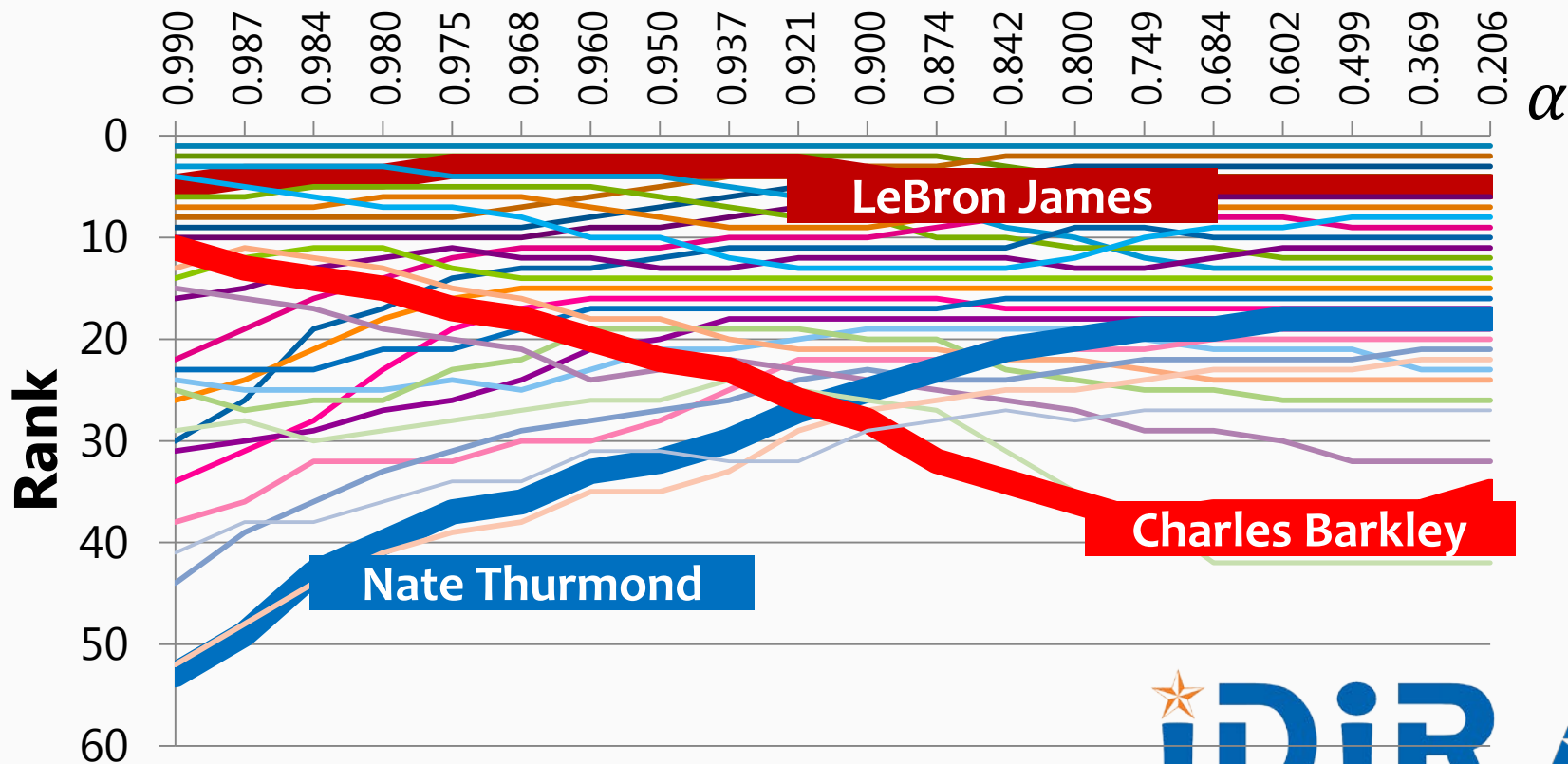
$\tau=6 \rightarrow$ 2-skyband (size 5; 3-skyband would be too big)



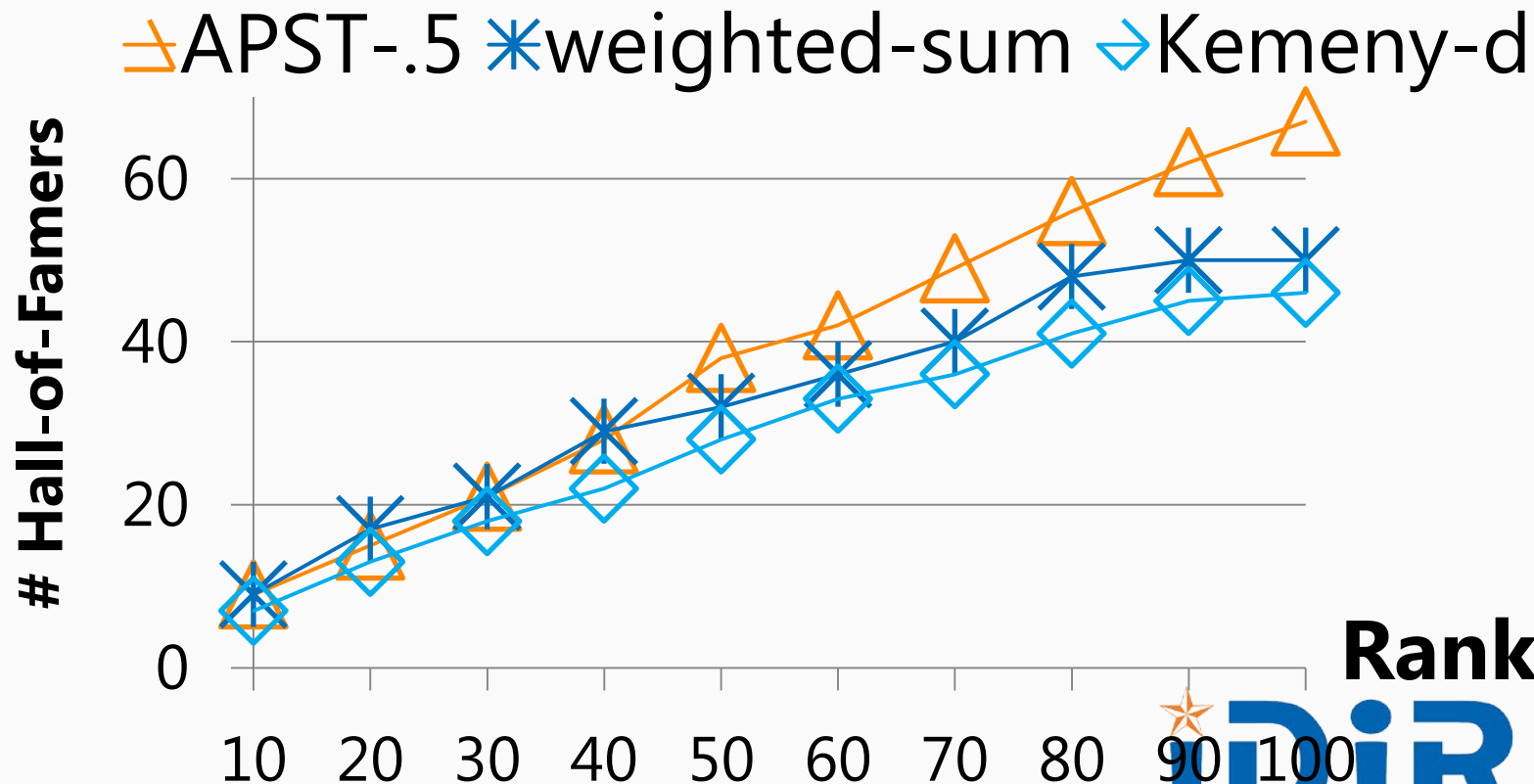
Experiments



Experiments



Experiments



Acknowledgment

UTA Students

- Naeemul Hassan
- Afroza Sultana
- Gensheng Zhang

Collaborators

- Pankaj Agarwal (Duke)
- Ping Luo (Chinese Academy of Sciences)
- Min Wang (Google Research)
- Jun Yang (Duke)
- Cong Yu (Google Research)



Acknowledgment

Funding sponsors



UNIVERSITY OF
TEXAS
ARLINGTON

Disclaimer: This material is based upon work partially supported by the National Science Foundation Grants 1018865, 1117369 and 1408928, 2011 and 2012 HP Labs Innovation Research Awards, and the National Natural Science Foundation of China Grant 61370019. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the funding agencies.



Thank You! Questions?

<http://ranger.uta.edu/~cli>

<http://idir.uta.edu>

cli@uta.edu

Demo

<http://idir.uta.edu/factwatcher>

