



**University of
Zurich^{UZH}**

Planning Ahead: Stream-Driven Linked-Data Access under Update-Budget Constraints

Shen Gao, Daniele Dell'Aglio, Soheila Dehghanzadeh,
Abraham Bernstein, Emanuele Della Valle, Alessandra Mileo



University of
Zurich^{UZH}



POLITECNICO
DI MILANO

Insight 

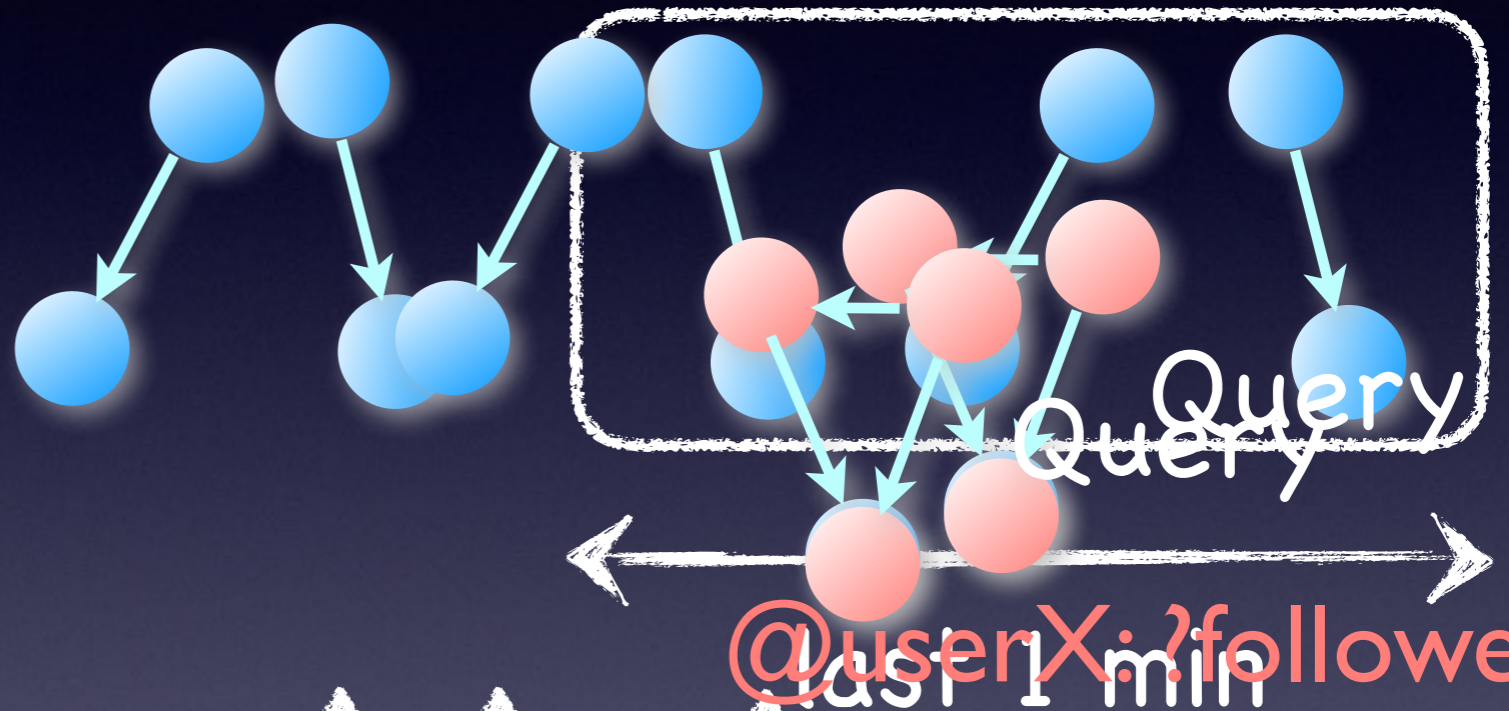
How can we efficiently
access SPARQL endpoints

in

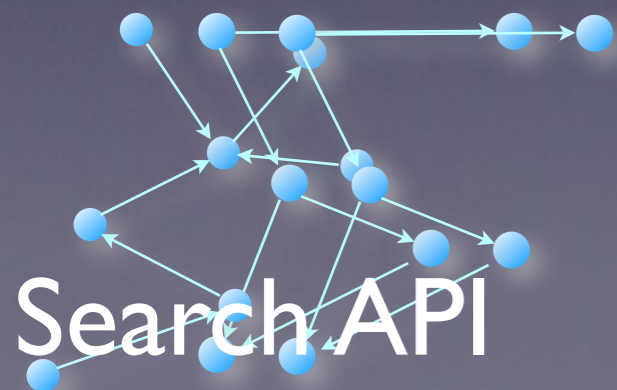
RDF *Stream* Processing (*RSP*)?

Finding Trendy Hashtags

@user1: #ISWC2016



Stream API



Search API

.....

Accessing Background Knowledge is Costly

Local Cache?

Budget: Num. of accesses per window

How can we allocate budget efficiently?

- Model the query: Bipartite Graph
- Exploiting the bipartite graph for solutions

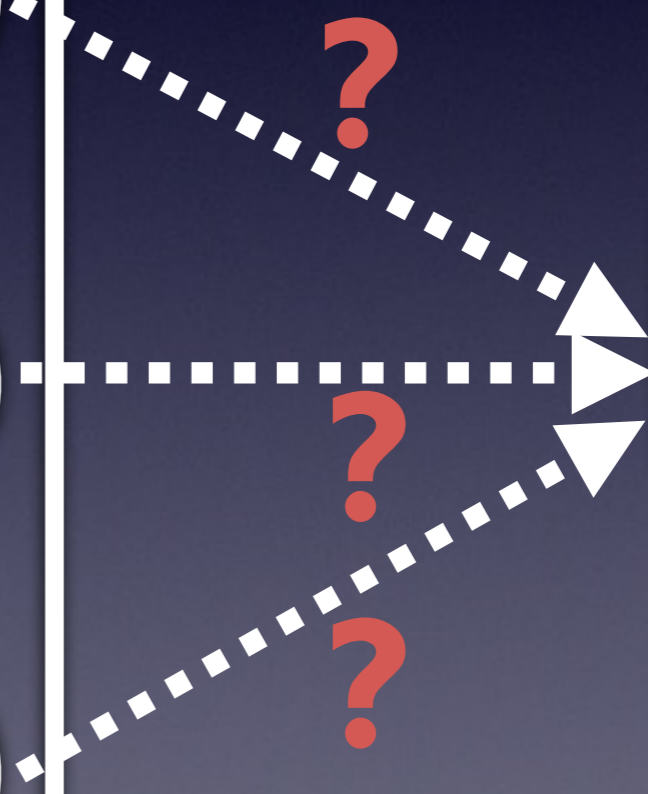
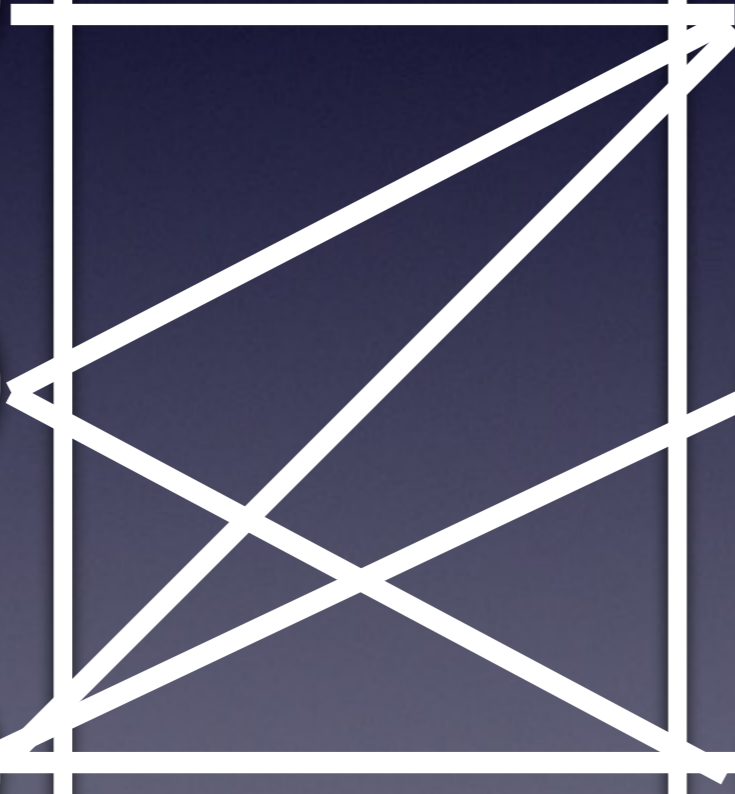
```
SELECT * WHERE  
        WINDOW (Stream) {PW}  
        SERVICE (BKG) {PS}
```

PW

PS

Stream

Local cache



Budget = 2

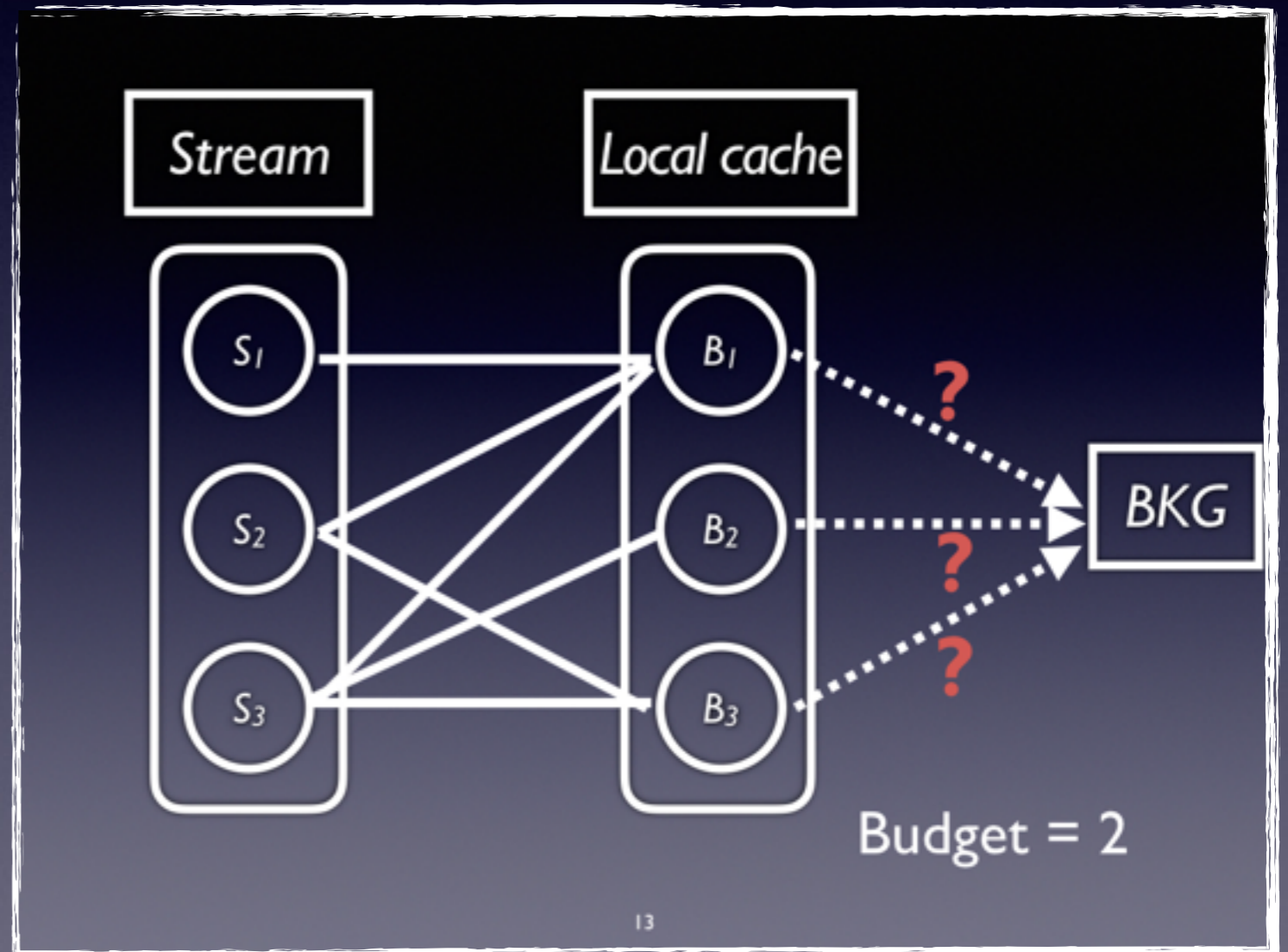
What is the optimisation goal?

```
SELECT * WHERE  
    WINDOW (Stream) {PW}  
    SERVICE (BKG) {PS}
```

- SERVICE clause P^S:
 - Basic Graph Pattern (BGP)
 - Aggregate query

Service Clause P^S - BGP

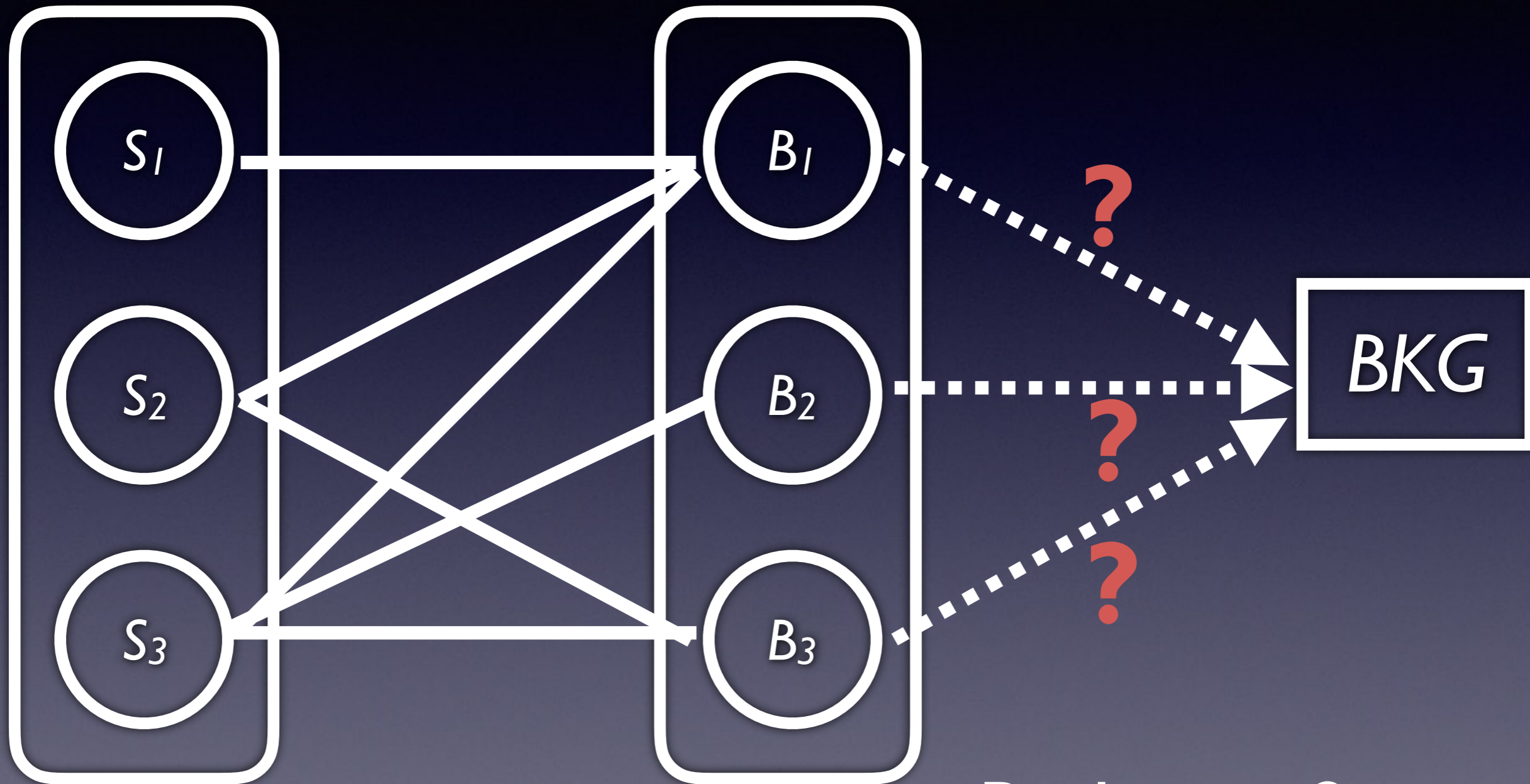
$\langle S_1 :p B_1 \rangle$
 $\langle S_2 :p B_1 \rangle$
 $\langle S_2 :p B_3 \rangle$
 $\langle S_3 :p B_1 \rangle$
 $\langle S_3 :p B_2 \rangle$
 $\langle S_3 :p B_3 \rangle$



Max(# of edges)

Stream

Local cache



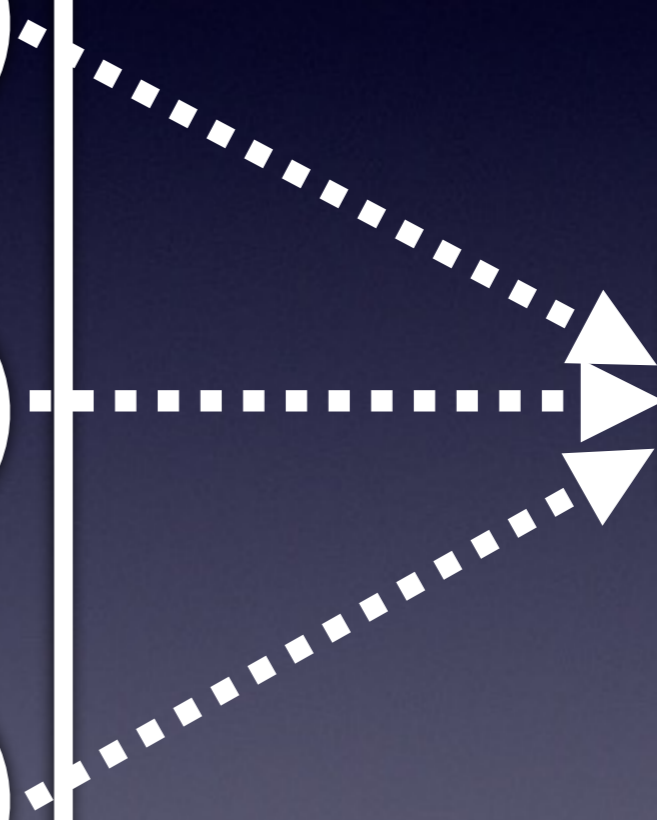
Budget = 2

Stream

Local cache



BKG

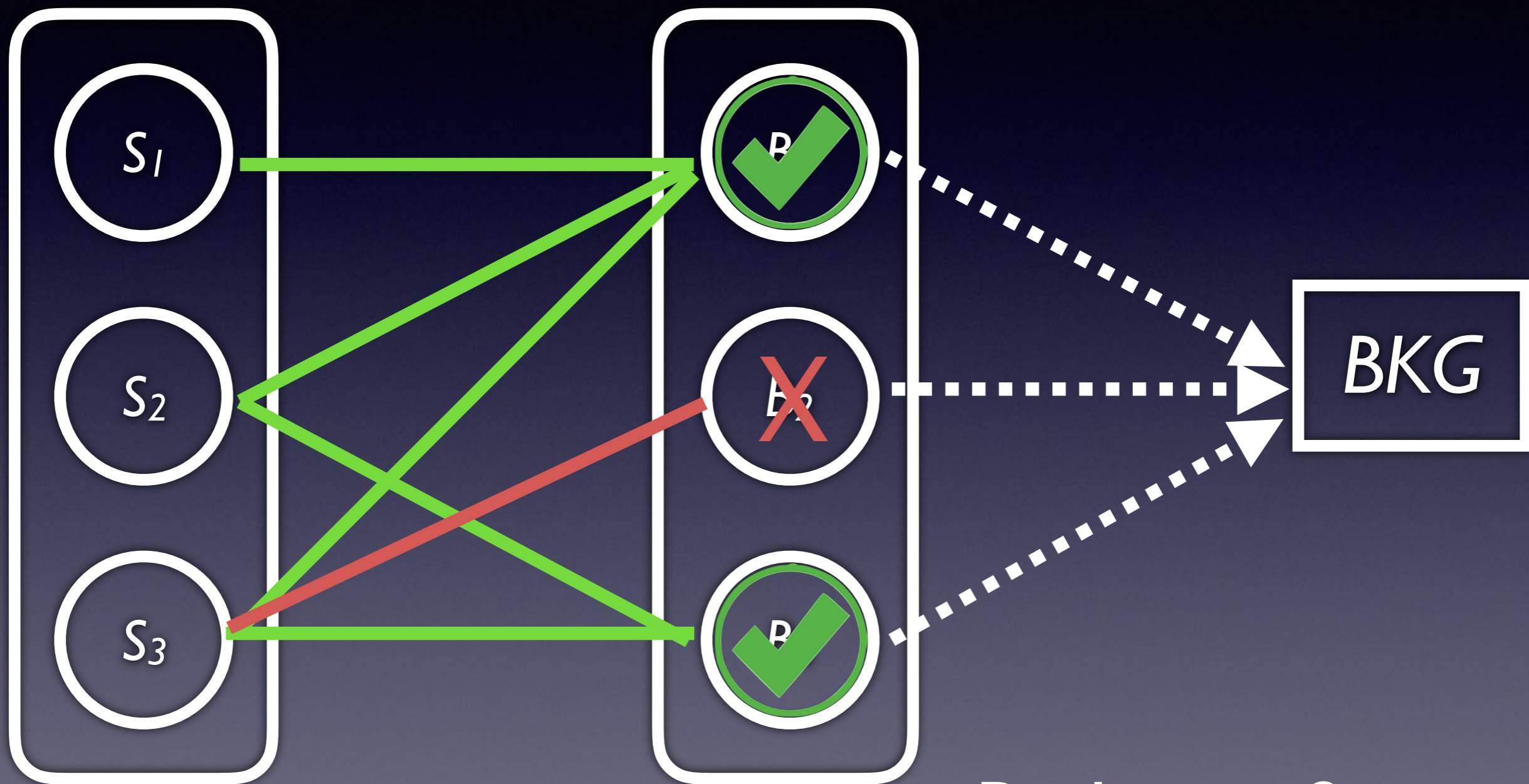


Budget = 2

of Fresh Results = 4

Stream

Local cache



Budget = 2

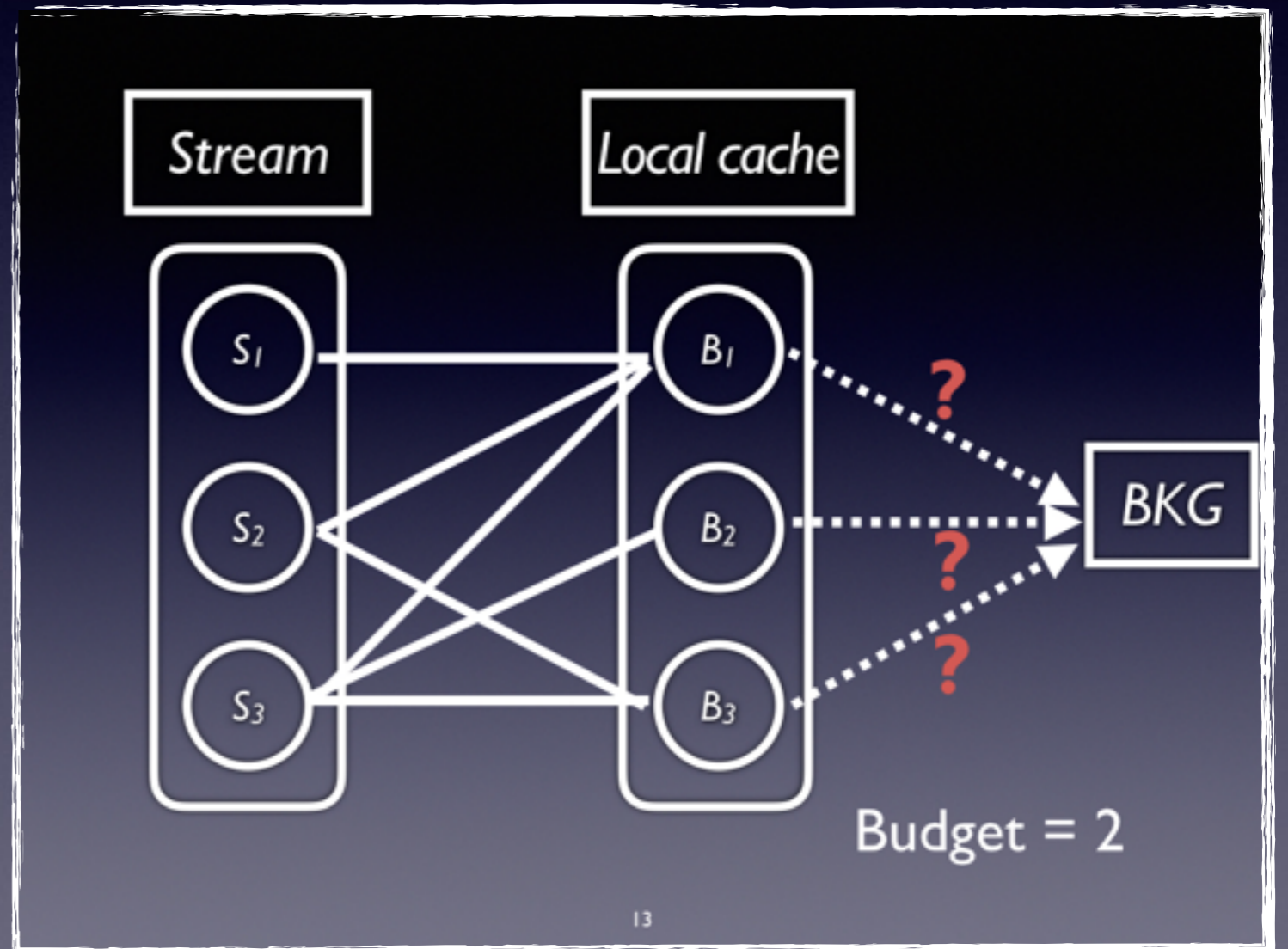
$SBM_{BGP} = \text{Max}(\text{deg}(B_i))$ # of Fresh Results = 5

Service Clause P^S - Agg.

$\langle S_1 : p \text{ Sum}(B_1) \rangle$

$\langle S_2 : p \text{ Sum}(B_1, B_3) \rangle$

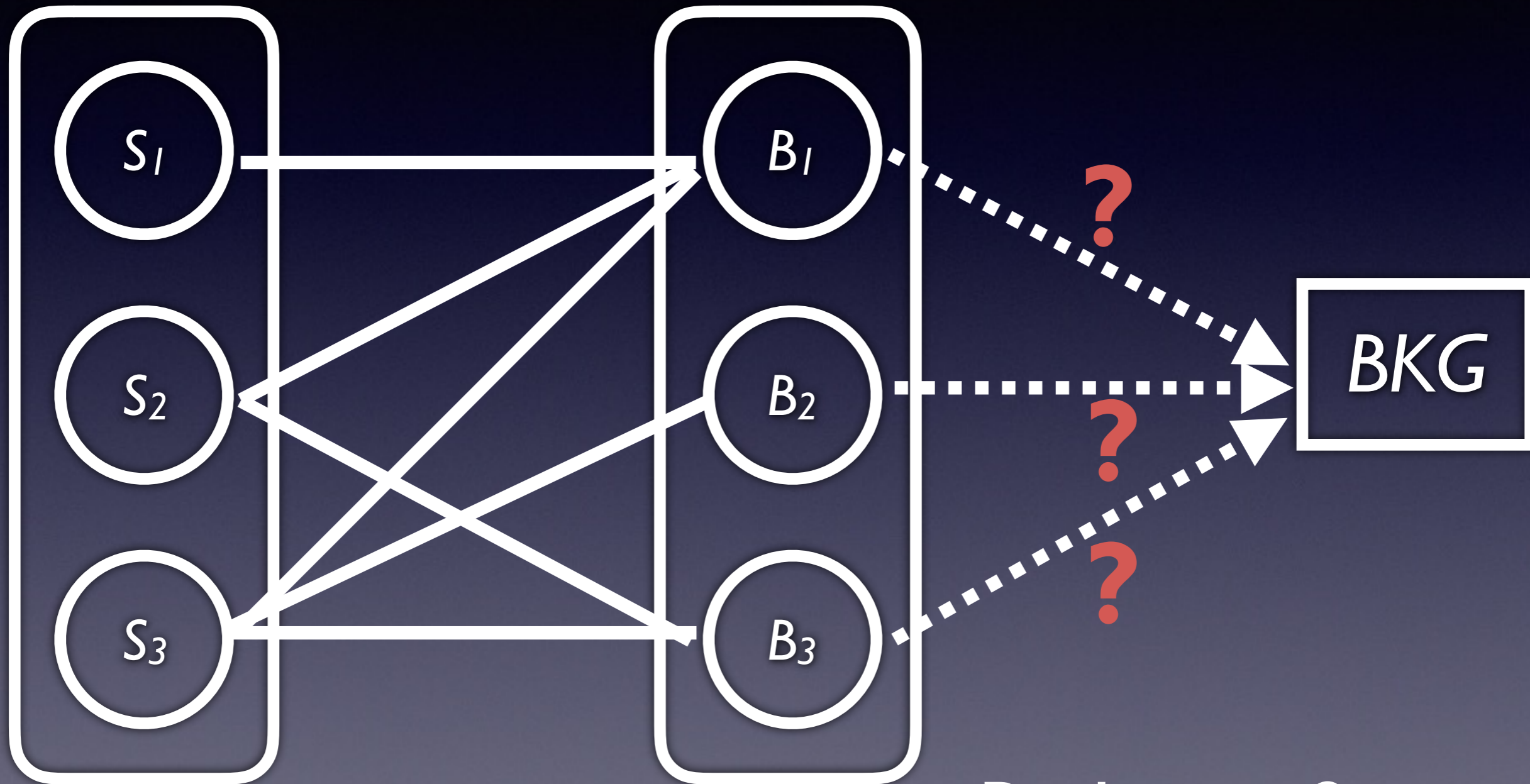
$\langle S_3 : p \text{ Sum}(B_1, B_2, B_3) \rangle$



Max(# of S_i , whose join partners are all fresh)

Stream

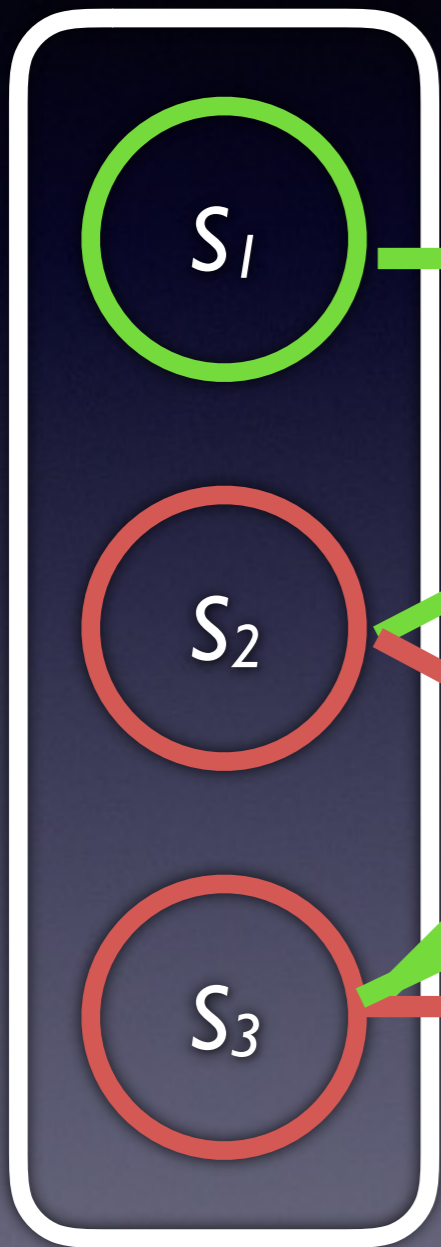
Local cache



Budget = 2

Stream

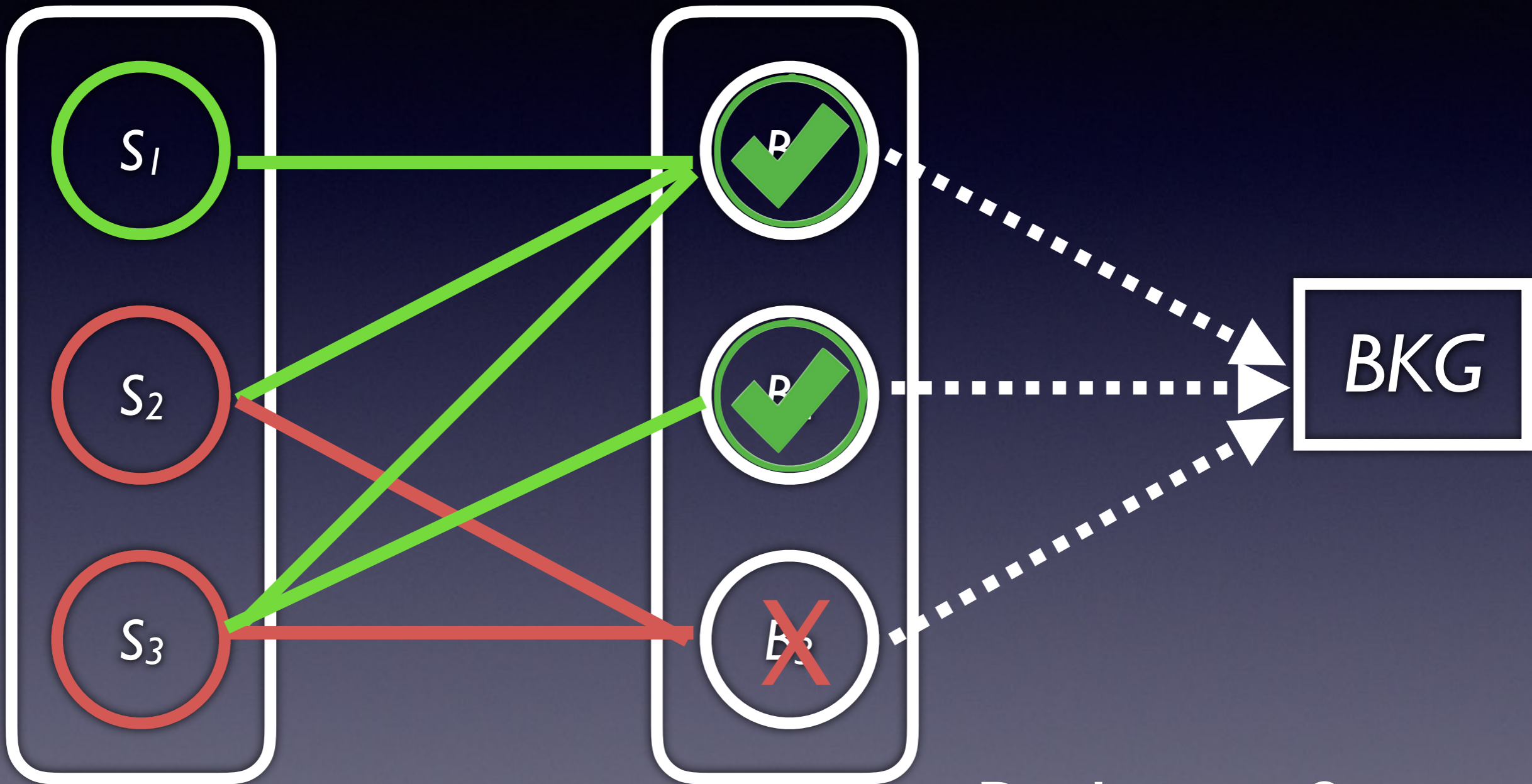
Local cache



BKG

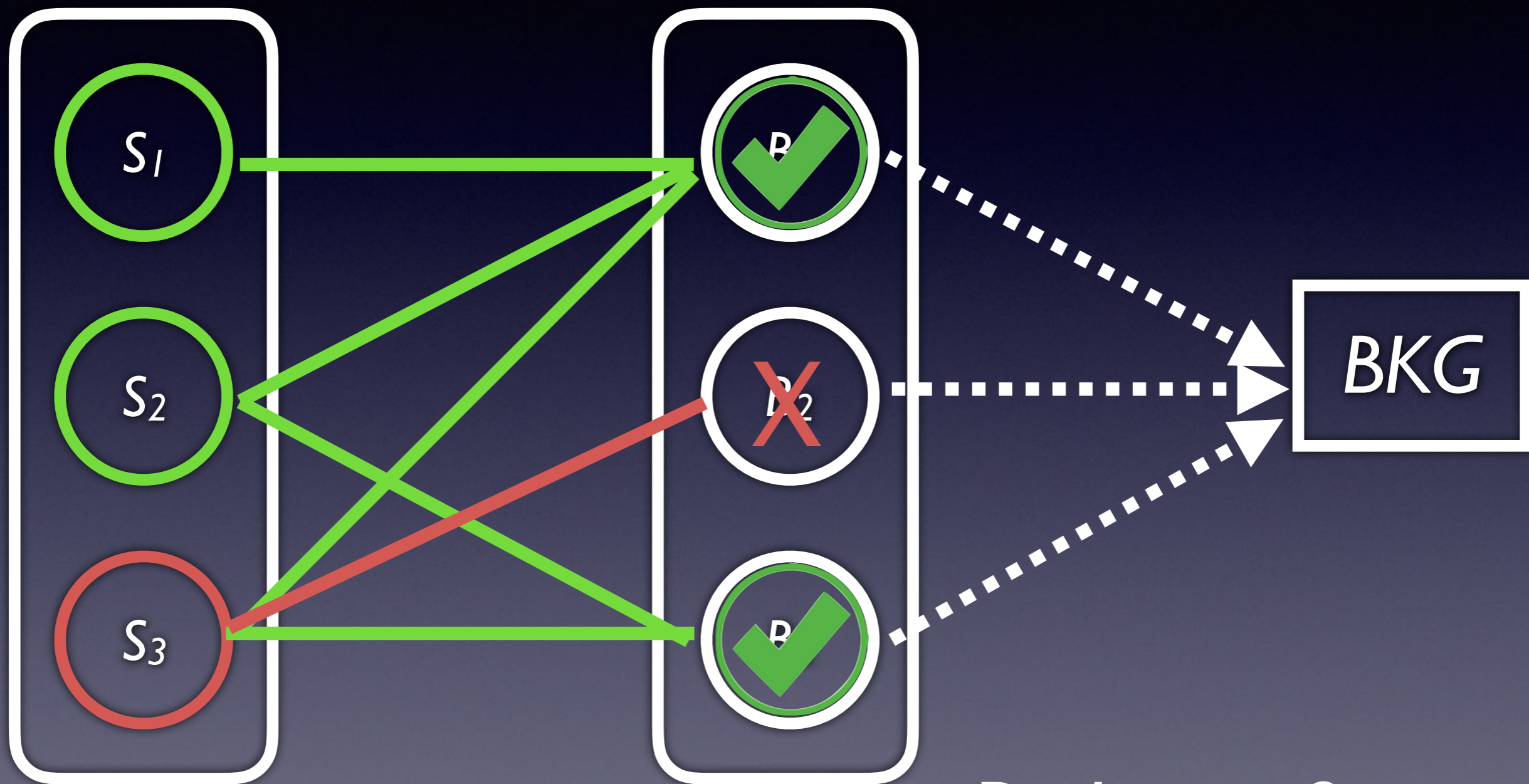
Budget = 2

of Fresh Results = 1



Stream

Local cache



NP-Hard

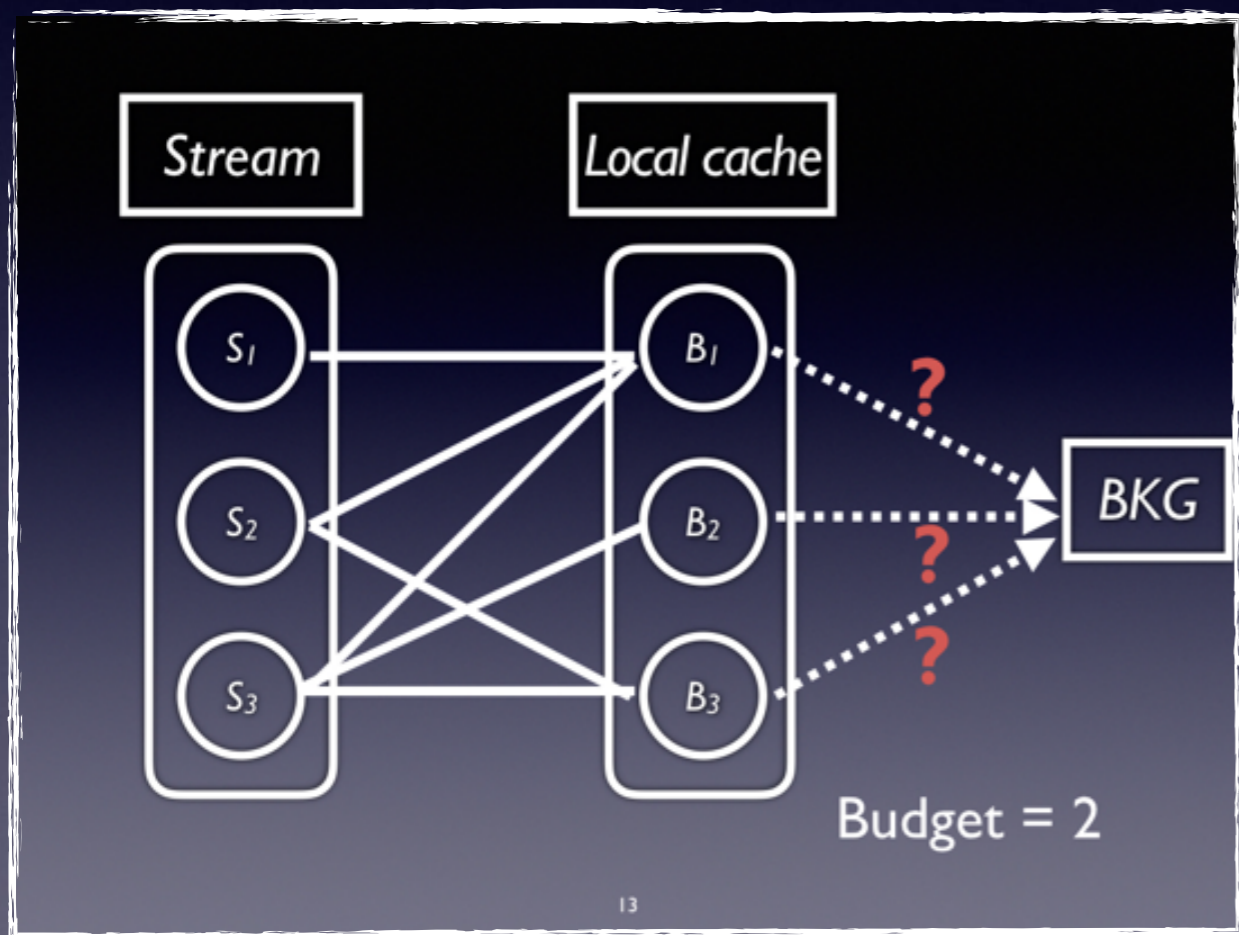
Budget = 2

$SBM_{agg} = B_i$ of $\text{Min}(\text{deg}(S_i))$ # of Fresh Results = 2

SBM:

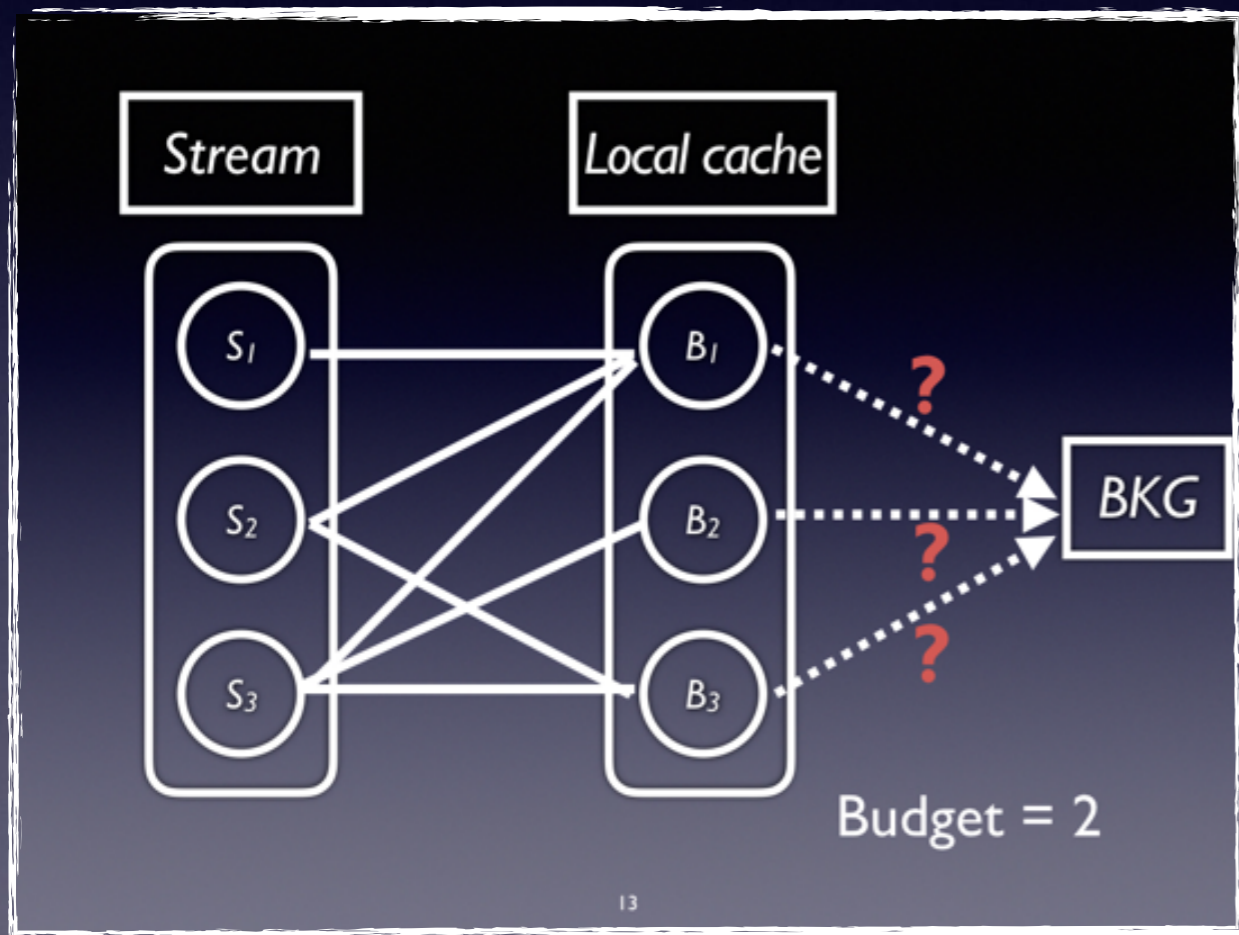
Selectivity Based Maintenance

BGP considers $\text{deg}(B_i)$
Agg. considers $\text{deg}(S_i)$
in the
current window



IBM:

Impact Based Maintenance



in the
current window

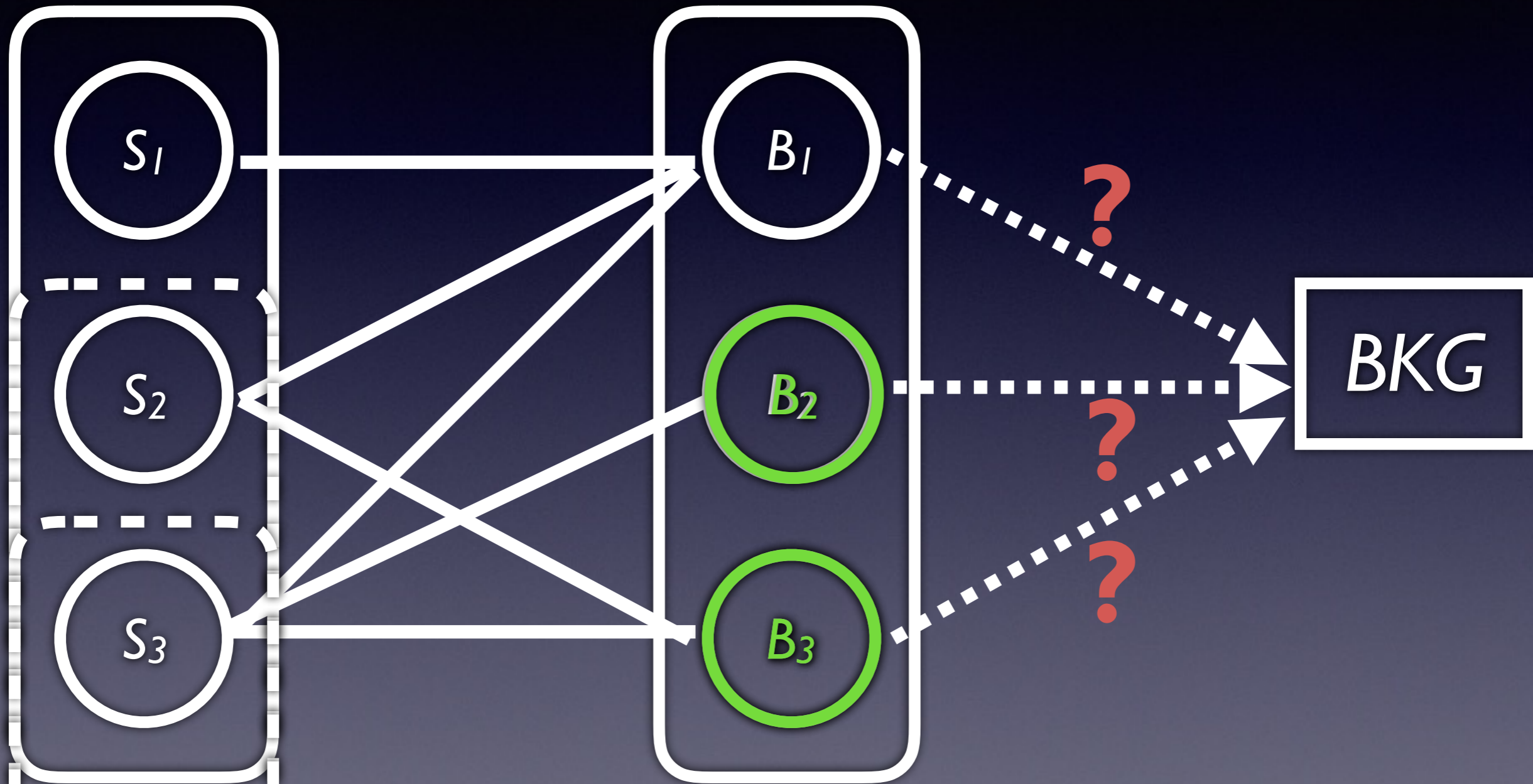
+

in the

following windows

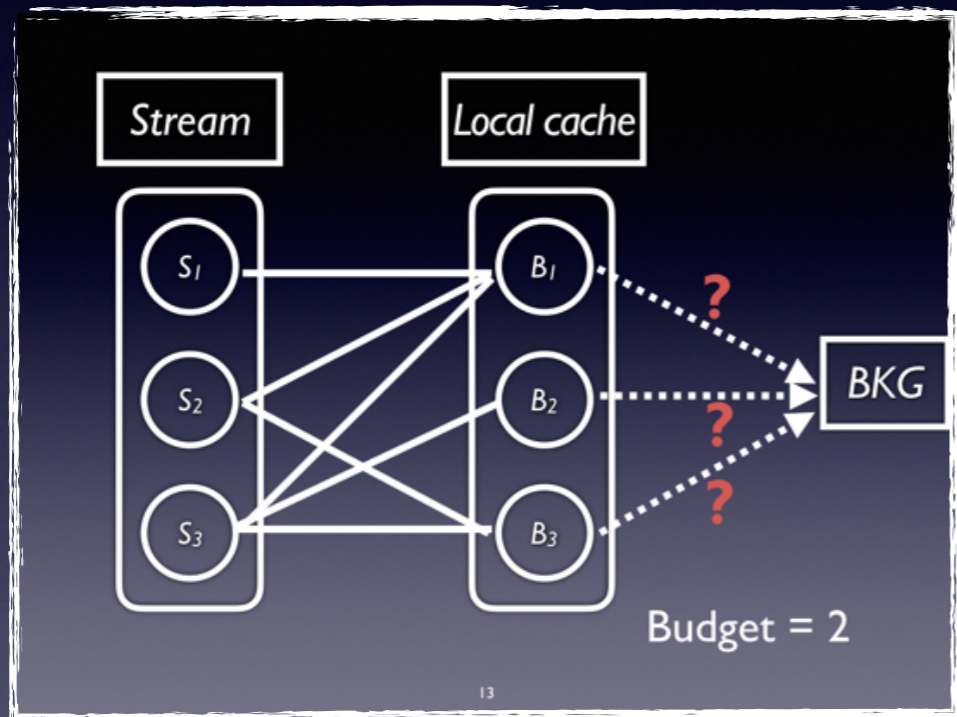
Stream

Local cache



Impact = # of results
in current + future windows

FBA: Flexible Budget Allocation



the budget

in the
current window
+
in the
following windows

Stream



Fixed Budget = 2

Fixed Budget = 2

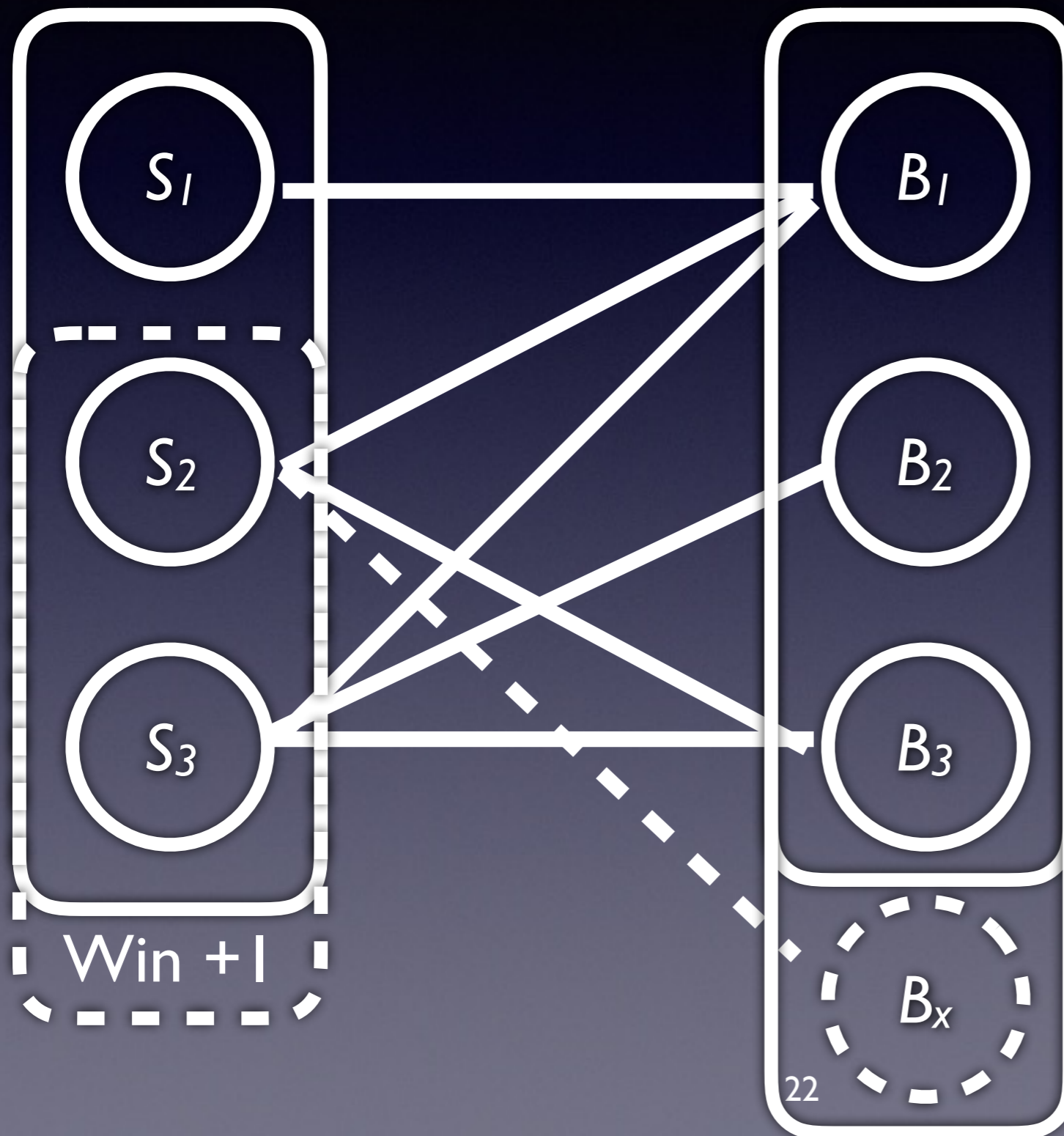
Fixed Budget = 2

Win +1

Win +2

Stream

Local cache



Stream



$$\text{Fixed Budget} = 2 - 1 = 1$$

$$\text{Fixed Budget} = 2 + 1 = 3$$

How much budget should be saved for future windows?

How can we allocate the budget better?

```
SELECT * WHERE  
        WINDOW (Stream) {PW}  
        SERVICE (BKG) {PS}
```

Spend the budget on
the ones that have
the largest impact over the windows!

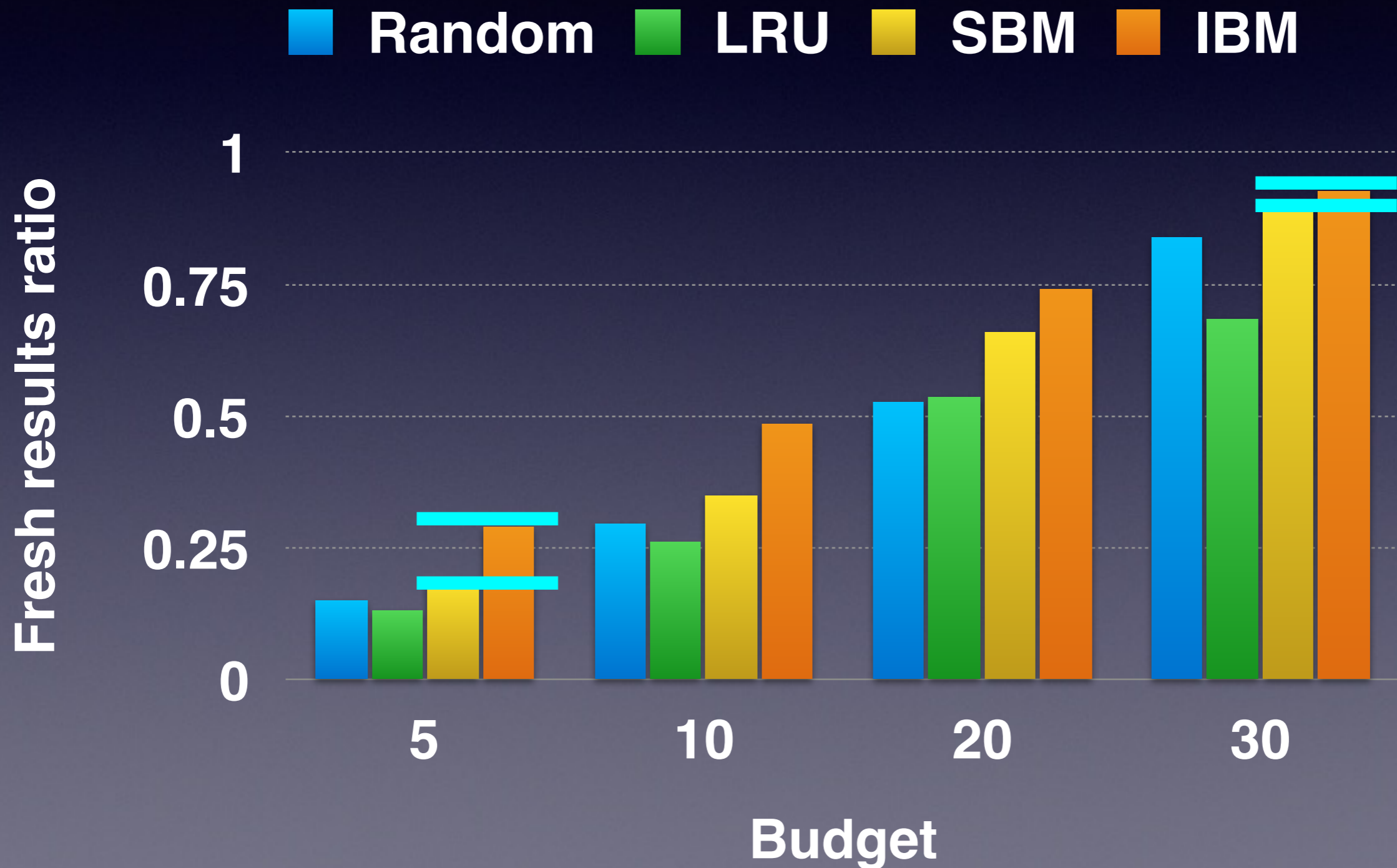
SBM => IBM => FBA

Experiment setup

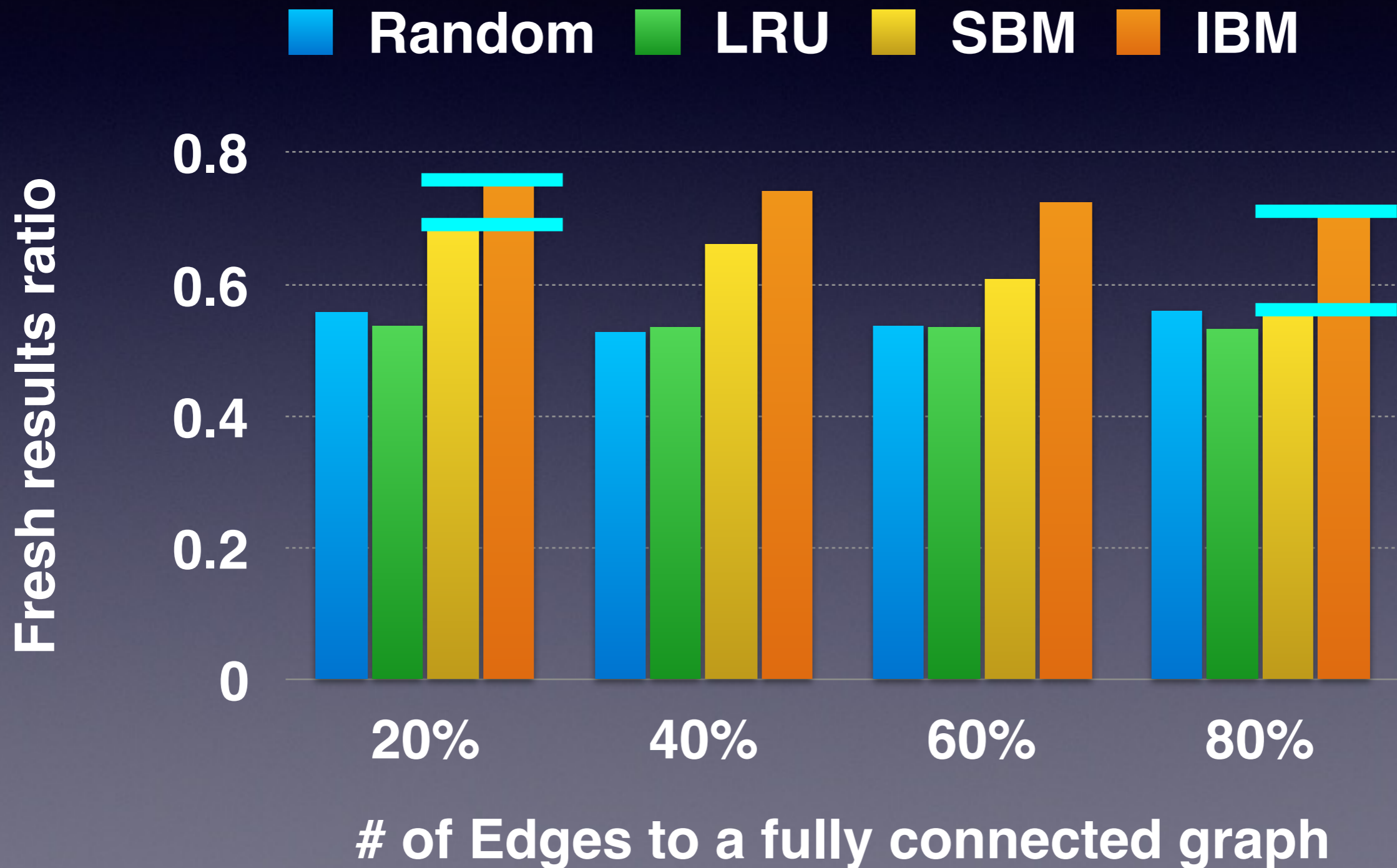
- Real implementation in C-SPARQL :
 - <https://github.com/dellaglio/CSPARQL-engine>

- Synthetic and Real data

Increasing Budget



Increasing Graph Density

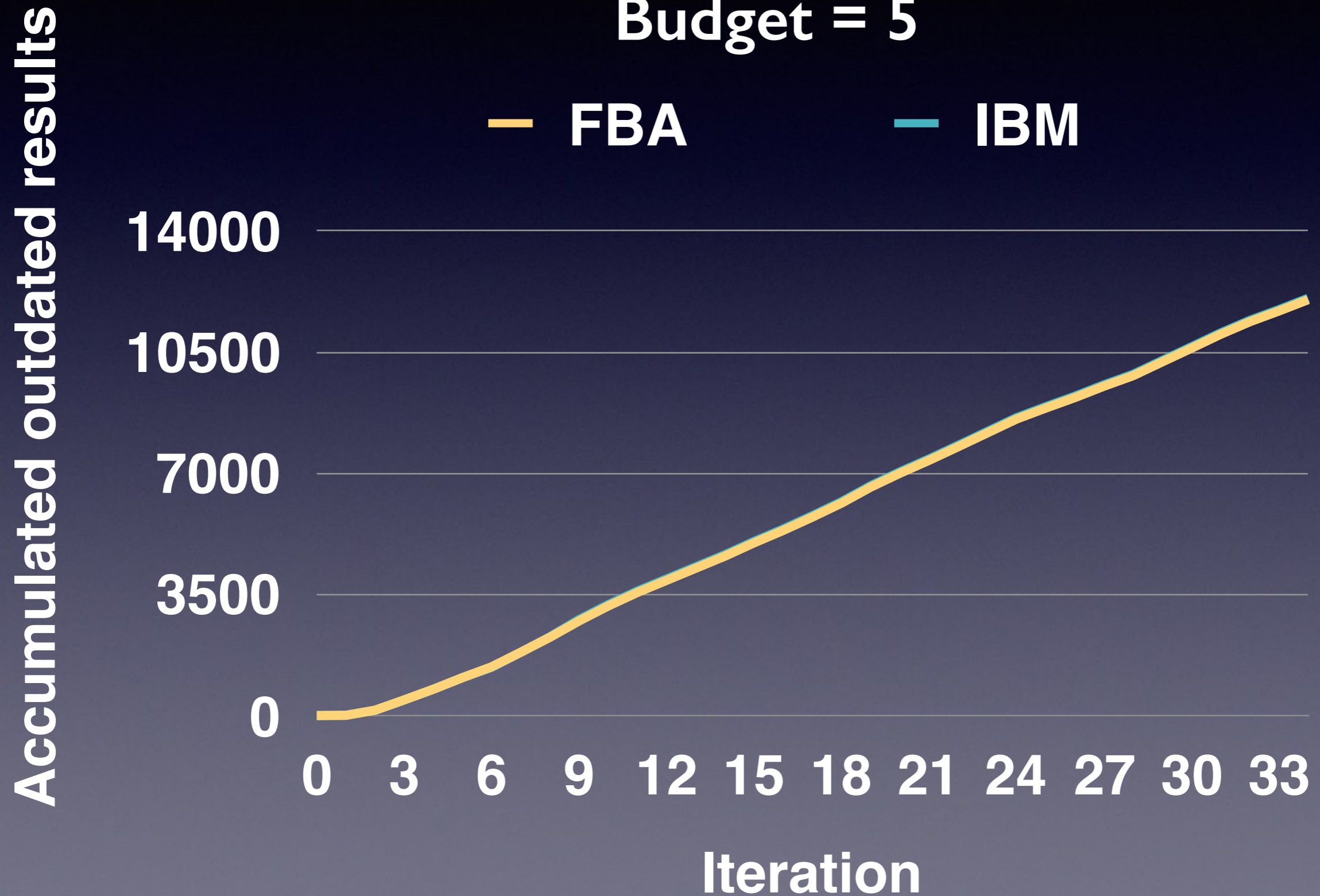


Increasing Budget

Budget = 5

— FBA

— IBM

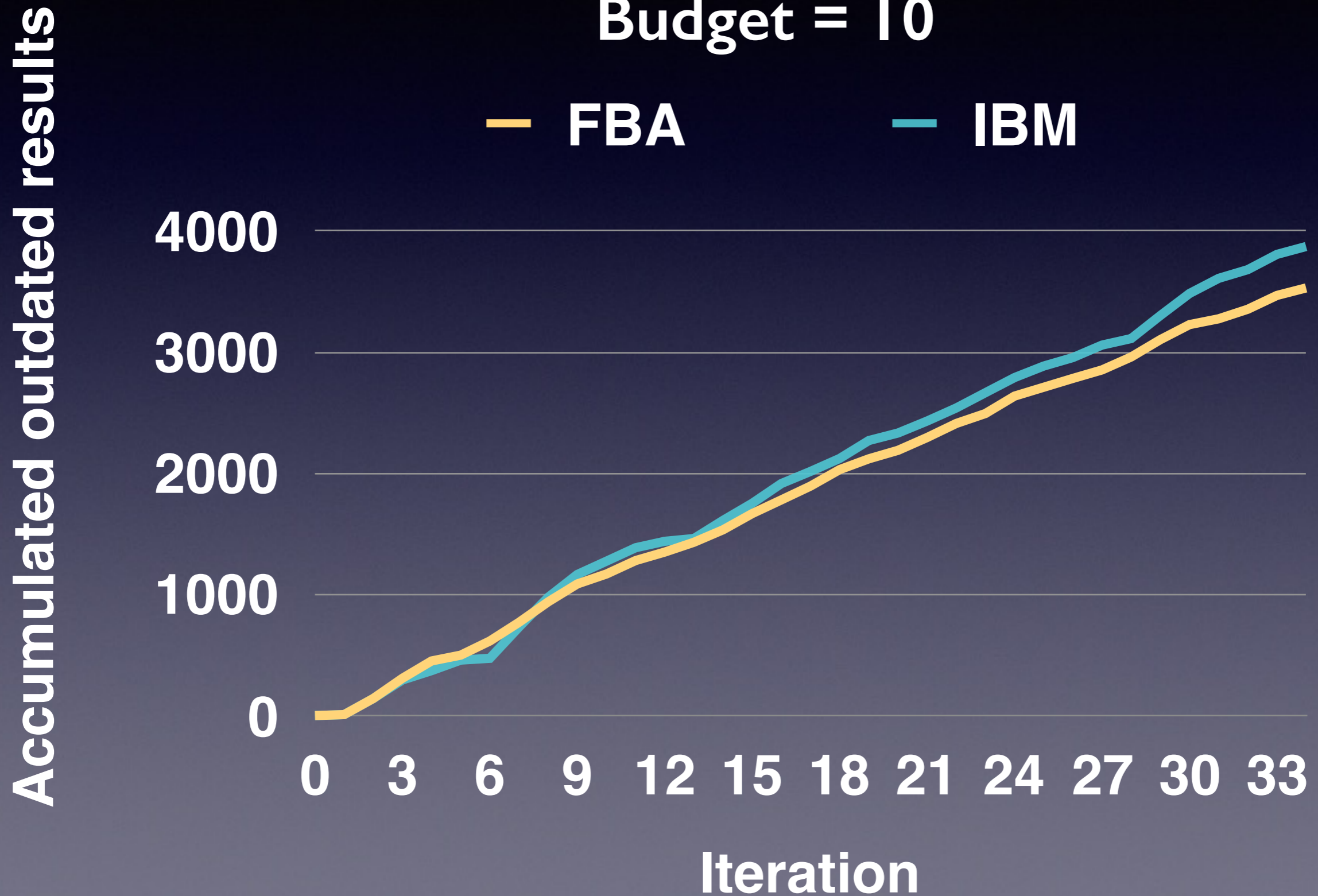


Increasing Budget

Budget = 10

— FBA

— IBM



Increasing Budget

Budget = 15

— FBA

— IBM

Accumulated outdated results

800

600

400

200

0

0

3

6

9

12

15

18

21

24

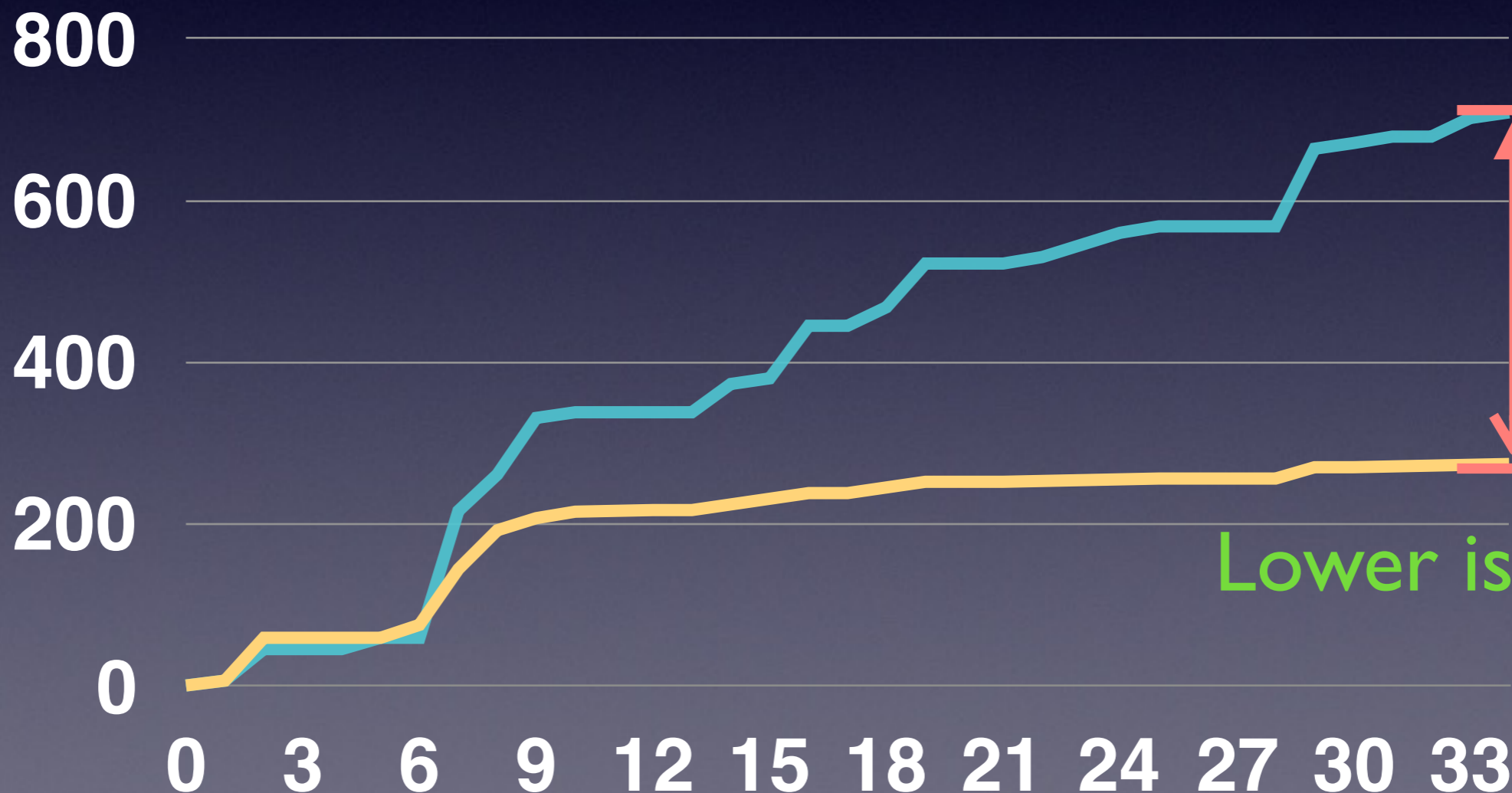
27

30

33

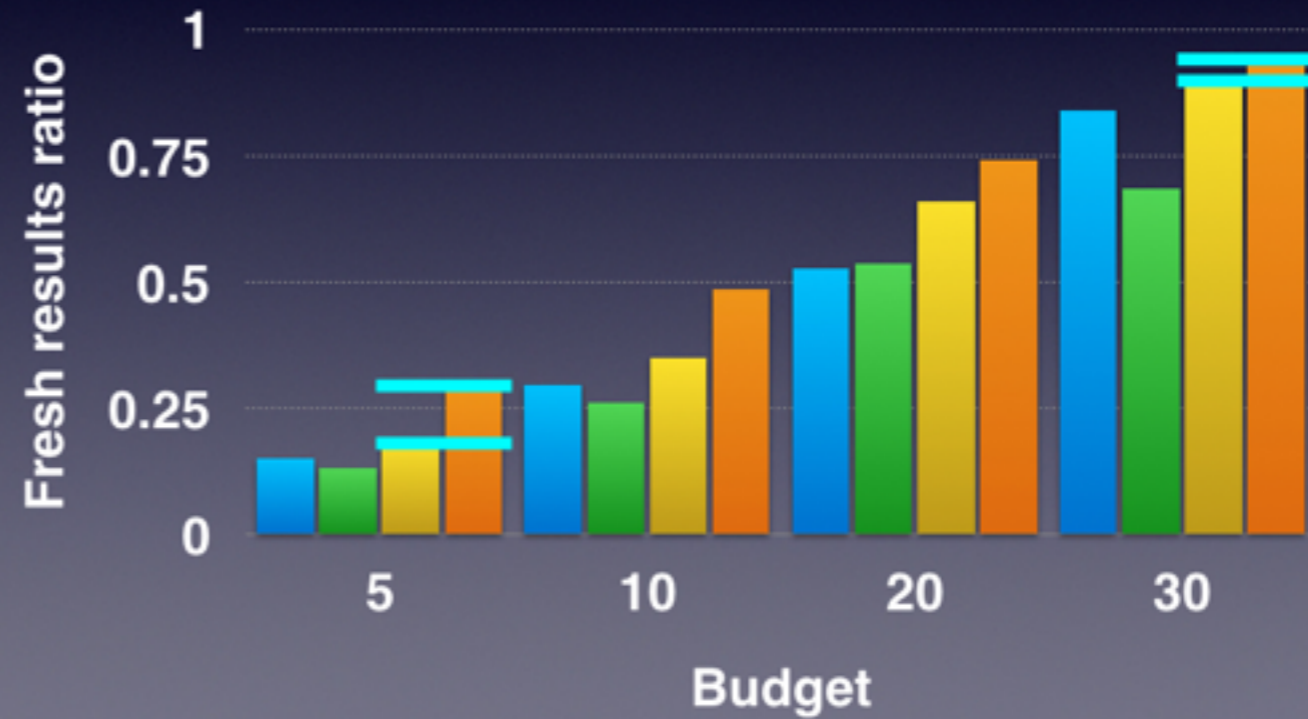
Iteration

Lower is better



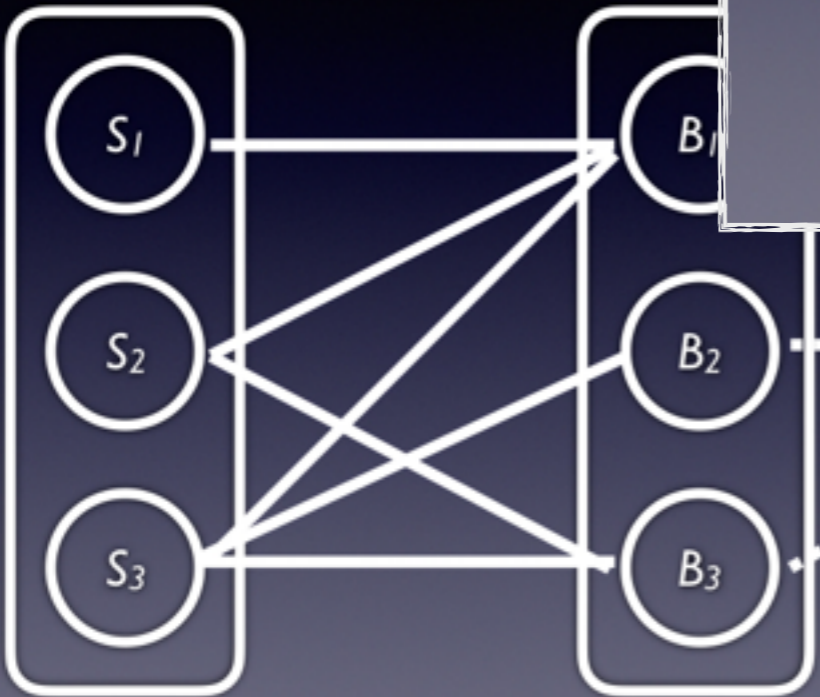
Increasing Budget

■ Random ■ LRU ■ SBM ■ IBM



Stream

Local d



BKG

Budget = 2

Thank you very much!

Questions?



**University of
Zurich^{UZH}**