# Uniform Split

## PostingFormat

Contributors
Bruno Roustant, Juan Camilo Rodriguez Duran, David Smiley
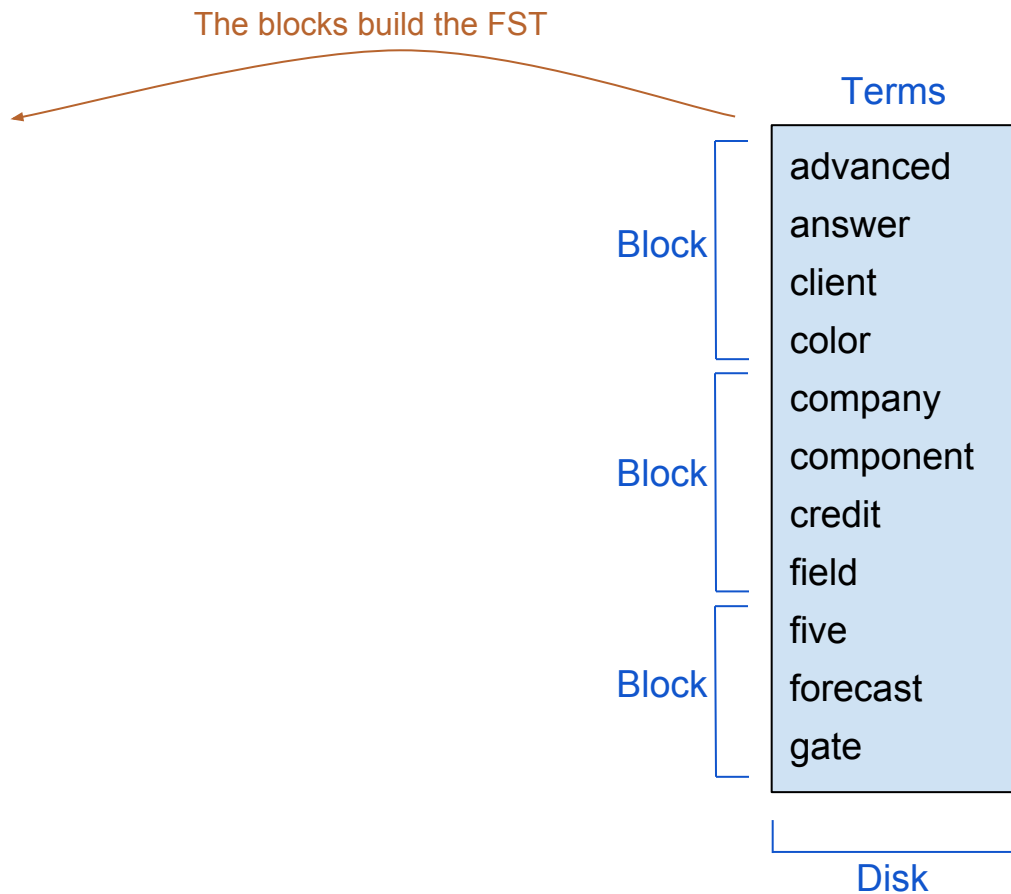
March, 2019

# Building the dictionary

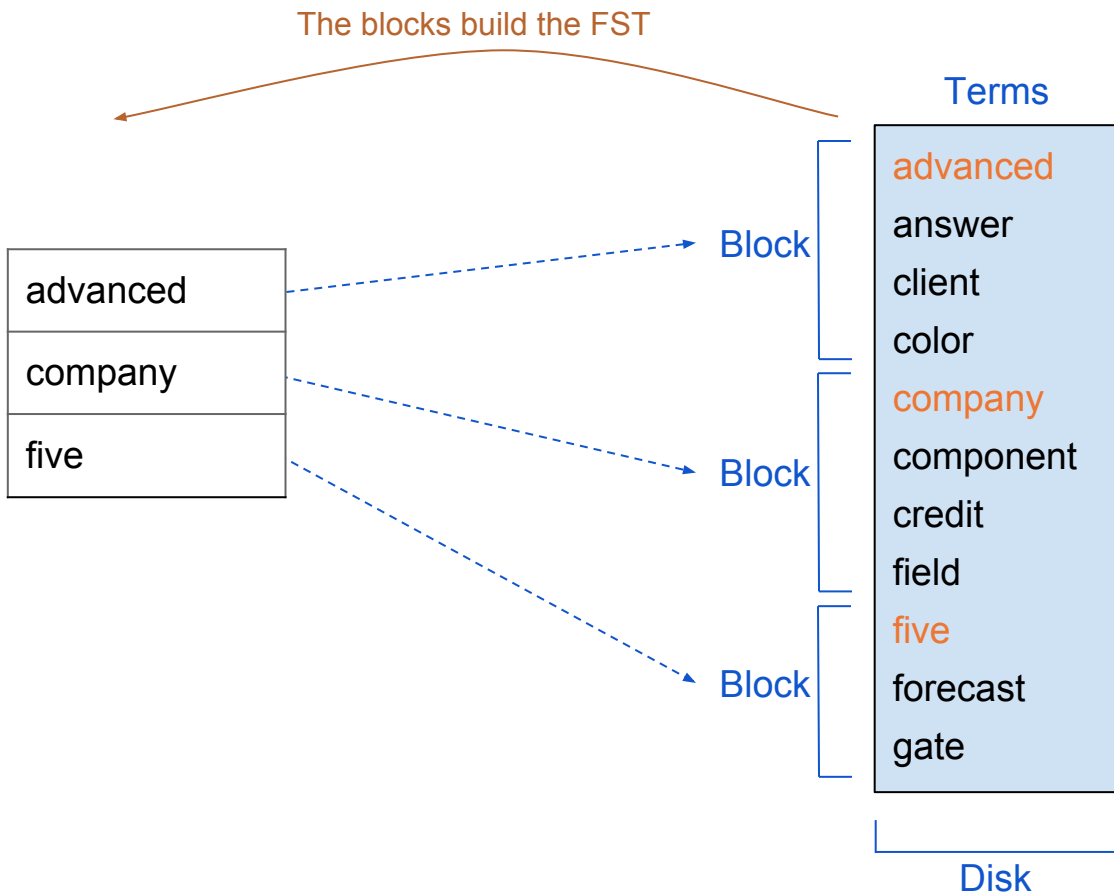The blocks of uniform size build the FST

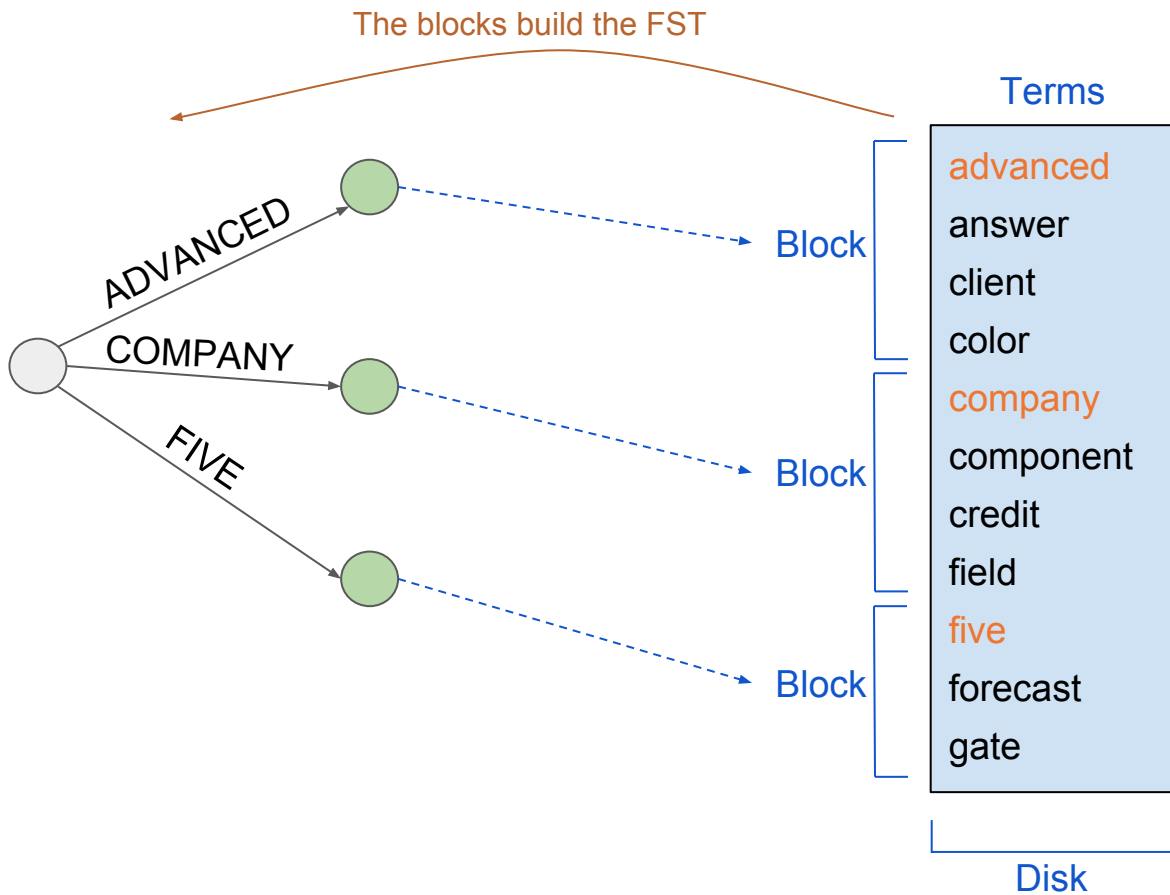Good size = 32

Size = 4 in this example

Terms

| Block | advanced |
| | answer |
| | client |
| | color |
| Block | company |
| | component |
| | credit |
| | field |
| Block | five |
| | forecast |
| | gate |

Disk

# Building the dictionary

The blocks of uniform size build the FST

The blocks build the FST

Terms

| advanced |
| company |
| five |

Block

Block

Block

advanced
answer
client
color
company
component
credit
field
five
forecast
gate

Disk

# Building the dictionary

The blocks of uniform size build the FST

The blocks build the FST

Terms

advanced
answer
client
color
company
component
credit
field
five
forecast
gate

Block

Block

Block

ADVANCED

COMPANY

FIVE

Disk

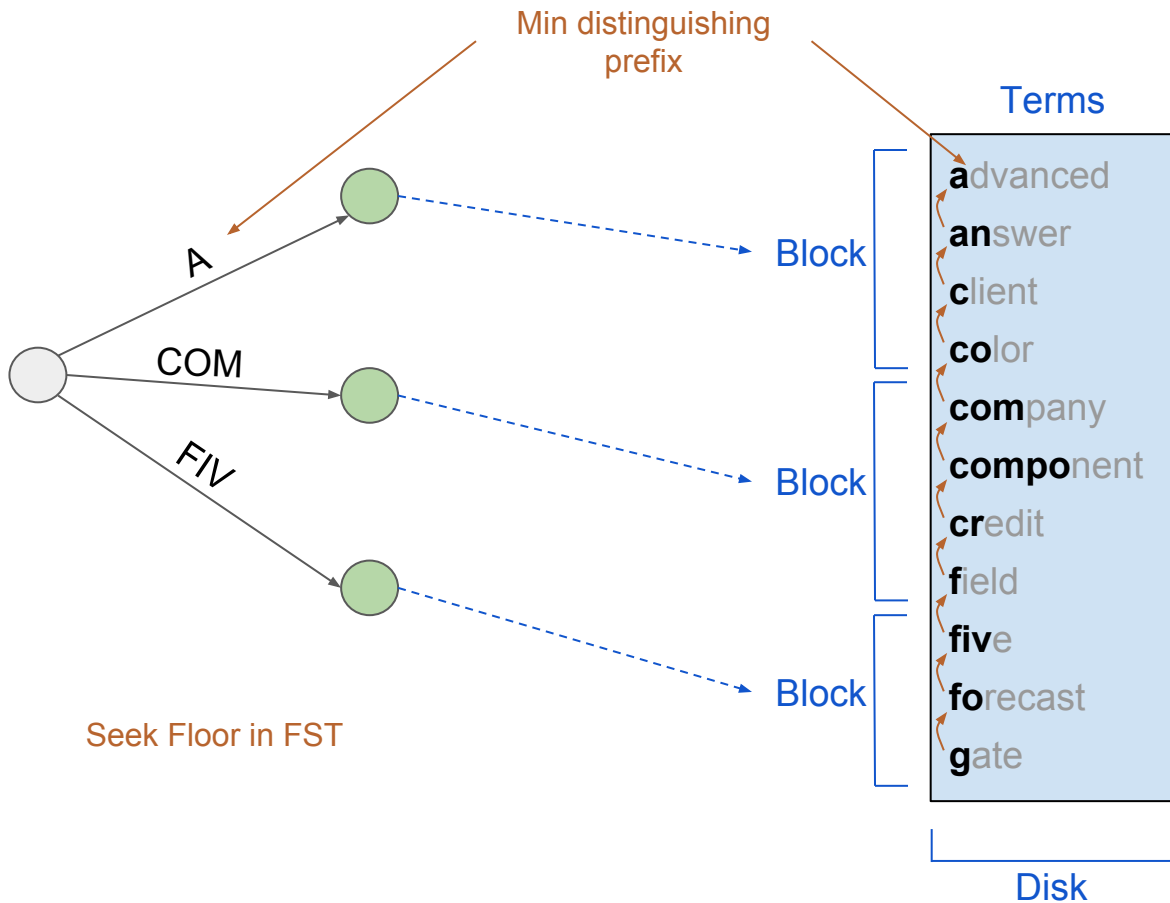# Building the dictionary

The blocks of uniform size build the FST

# Building the dictionary

## Optimization 1

Compute the **minimal distinguishing prefix (mdp)** for each term
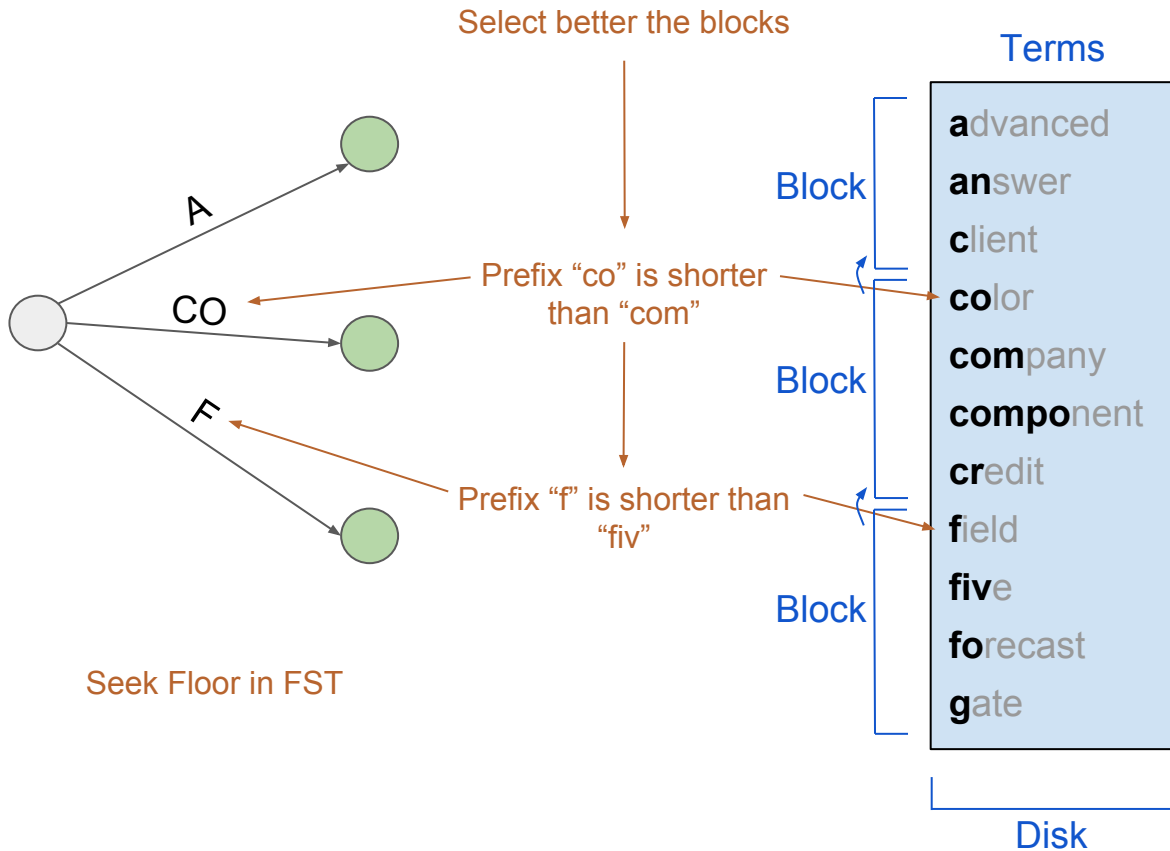
More compact
Same search

Min distinguishing prefix

Terms

Seek Floor in FST

**a**dvanced
**an**swer
**c**lient
**co**lor

**com**pany
**compo**nent
**cr**edit
**f**ield

**fiv**e
**fo**recast
**g**ate

A

COM

FIV

Block

Block

Block

Disk

# Building the dictionary

## Optimization 2

**Delta** select the best minimal distinguishing prefix (mdp) per block
Target block size +- 10%

## More compact
Same search

Select better the blocks

Terms

A

CO

Prefix "co" is shorter than "com"

F

Prefix "f" is shorter than "fiv"

Seek Floor in FST

Block

**a**dvanced
**an**swer
**c**lient
**co**lor
**com**pany
**compo**nent
**cr**edit
**f**ield
**fiv**e
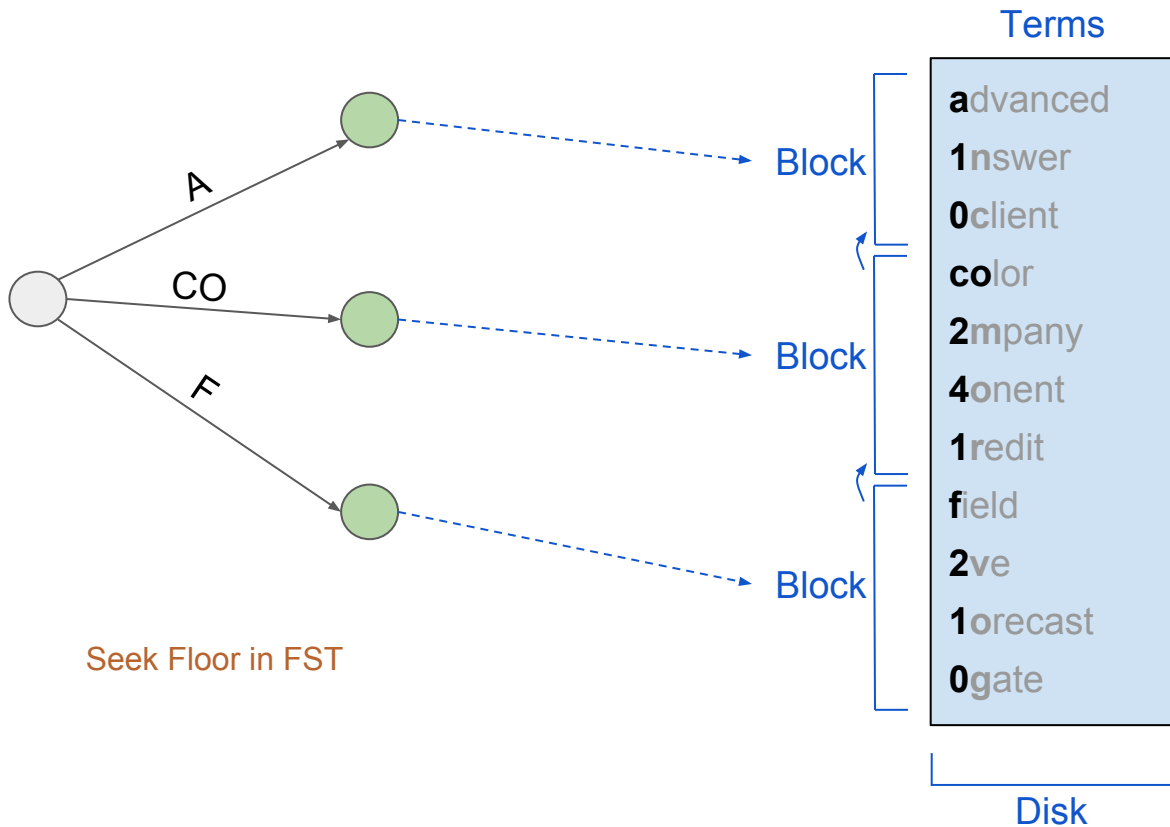**fo**recast
**g**ate

Block

Block

Disk

# Building the dictionary

## Optimization 3

**Incrementally** encode terms in each block
- Encode mdp length
- Only write suffix bytes

More compact
Faster scan

Seek Floor in FST

Terms

| | |
|---|---|
| Block | **a**dvanced |
| | **1**nswer |
| | **0**client |
| Block | **co**lor |
| | **2**mpany |
| | **4**onent |
| | **1**redit |
| Block | **f**ield |
| | **2**ve |
| | **1**orecast |
| | **0**gate |

Disk

# Building the dictionary

## Optimization 4

**Compare-&-jump to middle** term inside a block
- Cheaply reduces by half the block scan

## Faster scan

Seek Floor in FST

Terms

Block

Block

Block

**a**dvanced
**1n**swer
**0**client
**co**lor
**2m**pany
**4**onent
**1**redit
**f**ield
**2v**e
**1**orecast
**0**gate

Disk