# Data Transformations

data integration
and processing

data goes
through
operations

apply function
work from one format to another
join data with other datasets
filter values out of a data set

Transformations
* some are
aggregations

?

# After this video you will be able to..

- List common data transformations within big data pipelines

- Design a conceptual data processing pipeline using the basic data transformations
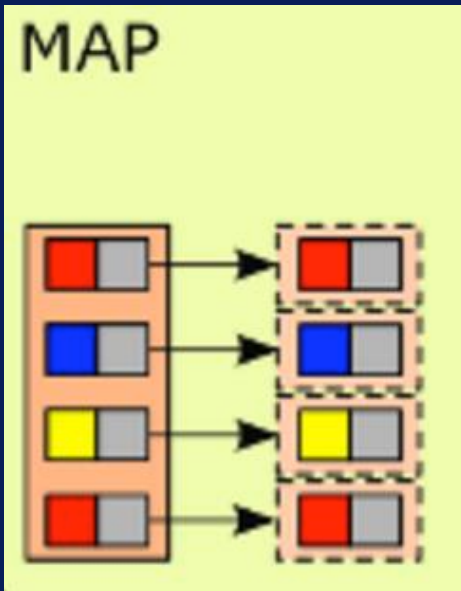
# Transformation are Tools to shape your data

*functions*

*convert from one form to another*

# Map — basic

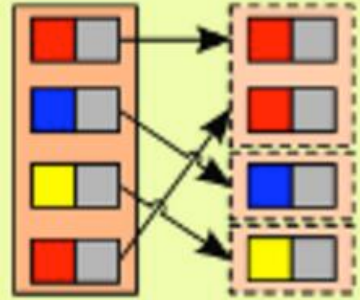

MAP

**Apply same operation to each member of a collection**

- **Color each member of a set**
- **Discount each product's price by 5%**
- **Apply formatting to each document in a folder**

— each data set is executed separately

same operation to each element

# Reduce


REDUCE

## 'Collecting' things that have same 'key'
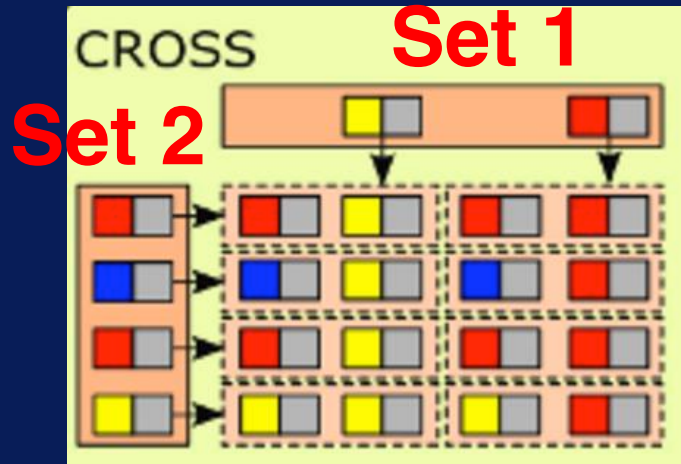
**Key: Colors**
- Collect blocks as per their colors

**Word Count Example**

**Key: Words**
- Sum frequency counts of words

— collectively apply the same process to objects of similar nature

* map and reduce are types of transformations that work on a single list of key and data pairings
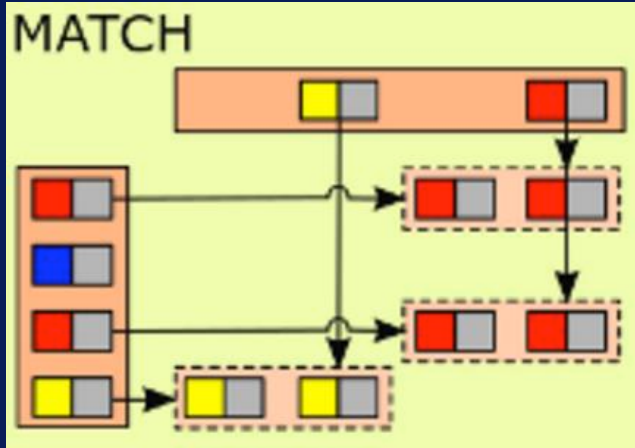
# Cross/ Cartesian



→ Multiplication

**Do some process to each pair from two sets**

# Match/Join



→ Selective Multiplication

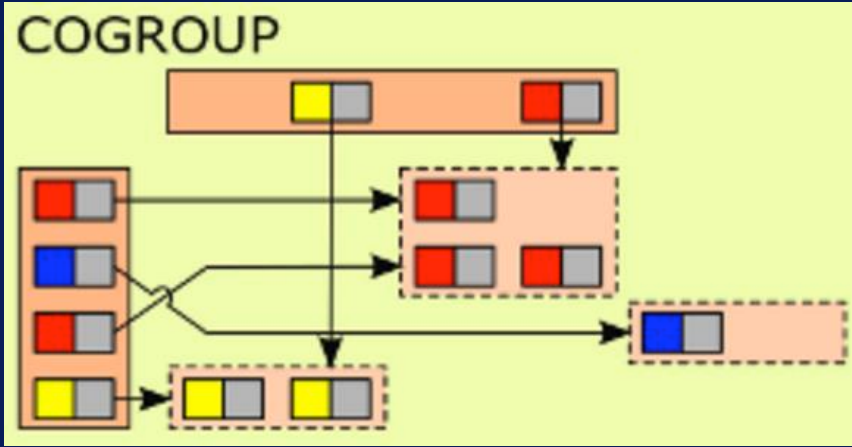**Do some process to each <u>pair</u> from two sets – which have same 'key'**

*just grouping together the data partitions with the same key*

*\* selective in forming pairs — every pair must have something in common Key*
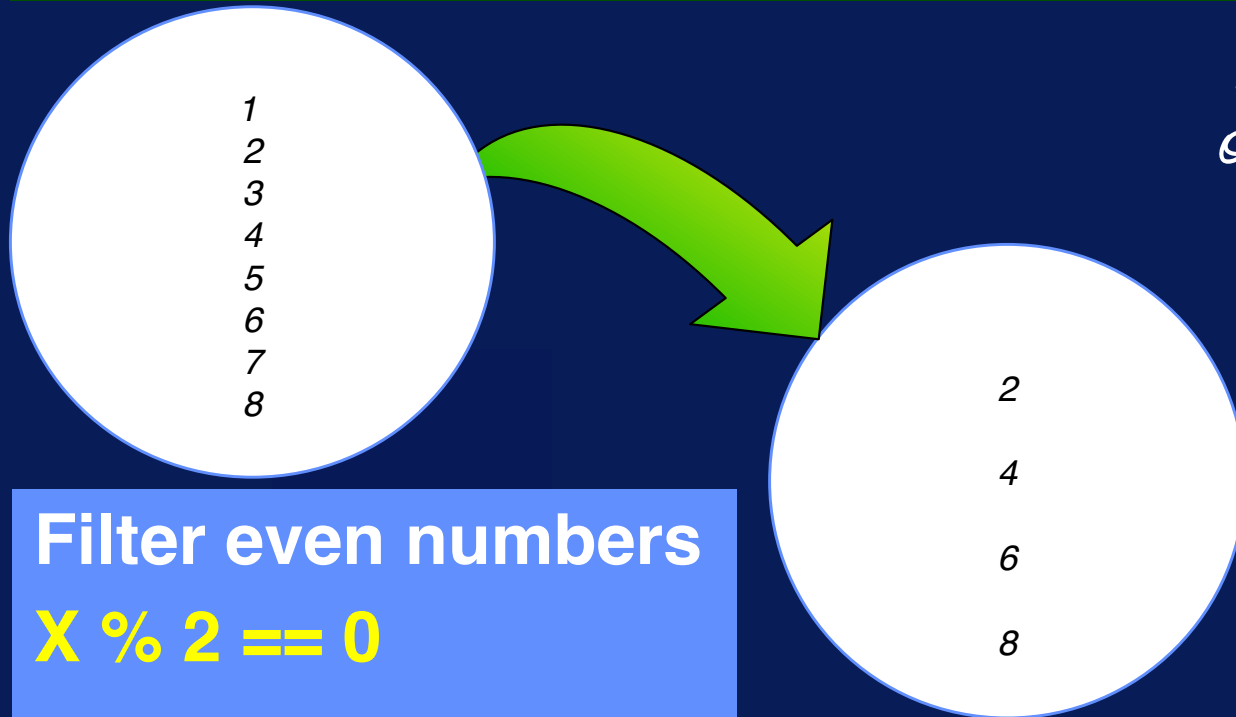
# Co-Group



→ Group common items

- **Collect similar things first**
- **Apply a process to each collection**
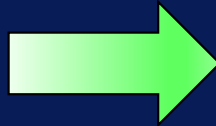
*listing __even__ if they dont exist in both datasets*

# Filter

Select elements that match a criteria

1
2
3
4
5
6
7
8

Filter even numbers
X % 2 == 0

2
4
6
8

like a test,
only elements that
pass the-test
are shown

# Basic Transformations → Get Results



effectiveness of transformation
is in pipelining them in a
way that helps solve problem
as you would perform a series
of tasks on a real block