# 61A Lecture 36

Wednesday, November 30

# Project 4 Contest Gallery

Tuesday, November 29, 2011

# Project 4 Contest Gallery

Prizes will be awarded for the winning entry in each of the
following categories.

# Project 4 Contest Gallery

Prizes will be awarded for the winning entry in each of the following categories.

- **Featherweight.** At most 128 words of Logo, not including comments and delimiters.

Tuesday, November 29, 2011

# Project 4 Contest Gallery

Prizes will be awarded for the winning entry in each of the following categories.

- **Featherweight.** At most 128 words of Logo, not including comments and delimiters.

- **Heavyweight.** At most 1024 words of Logo, not including comments and delimiters.

# Project 4 Contest Gallery

Prizes will be awarded for the winning entry in each of the following categories.

- **Featherweight.** At most 128 words of Logo, not including comments and delimiters.

- **Heavyweight.** At most 1024 words of Logo, not including comments and delimiters.

Winners will be selected by popular vote! (Homework 13)

# Project 4 Contest Gallery

Prizes will be awarded for the winning entry in each of the following categories.

- **Featherweight.** At most 128 words of Logo, not including comments and delimiters.

- **Heavyweight.** At most 1024 words of Logo, not including comments and delimiters.

Winners will be selected by popular vote! (Homework 13)

- Static **images** of the output of your programs

# Project 4 Contest Gallery

Prizes will be awarded for the winning entry in each of the following categories.

- **Featherweight.** At most 128 words of Logo, not including comments and delimiters.

- **Heavyweight.** At most 1024 words of Logo, not including comments and delimiters.

Winners will be selected by popular vote! (Homework 13)

- Static **images** of the output of your programs

- Tonight at midnight: I'll post your Logo **implementations**!

Tuesday, November 29, 2011

# Project 4 Contest Gallery

Prizes will be awarded for the winning entry in each of the following categories.

- **Featherweight.** At most 128 words of Logo, not including comments and delimiters.

- **Heavyweight.** At most 1024 words of Logo, not including comments and delimiters.

Winners will be selected by popular vote! (Homework 13)

- Static **images** of the output of your programs

- Tonight at midnight: I'll post your Logo **implementations**!
  - Run them to see these images evolve!

# Project 4 Contest Gallery

Prizes will be awarded for the winning entry in each of the following categories.

- **Featherweight.** At most 128 words of Logo, not including comments and delimiters.

- **Heavyweight.** At most 1024 words of Logo, not including comments and delimiters.

Winners will be selected by popular vote! (Homework 13)

- Static **images** of the output of your programs

- Tonight at midnight: I'll post your Logo **implementations**!
  - Run them to see these images evolve!

- I will also post a **solution** to the Logo project

# Project 4 Contest Gallery

Prizes will be awarded for the winning entry in each of the following categories.

- **Featherweight.** At most 128 words of Logo, not including comments and delimiters.

- **Heavyweight.** At most 1024 words of Logo, not including comments and delimiters.

Winners will be selected by popular vote! (Homework 13)

- Static **images** of the output of your programs

- Tonight at midnight: I'll post your Logo **implementations**!
  - Run them to see these images evolve!

- I will also post a **solution** to the Logo project
  - It runs (almost) all of the contest entries

# Project 4 Contest Gallery

Prizes will be awarded for the winning entry in each of the following categories.

- **Featherweight.** At most 128 words of Logo, not including comments and delimiters.

- **Heavyweight.** At most 1024 words of Logo, not including comments and delimiters.

Winners will be selected by popular vote! (Homework 13)

- Static **images** of the output of your programs

- Tonight at midnight: I'll post your Logo **implementations**!
  - Run them to see these images evolve!

- I will also post a **solution** to the Logo project
  - It runs (almost) all of the contest entries
  - You can use it as a study guide for the final

Tuesday, November 29, 2011

# Project 4 Contest Gallery

Prizes will be awarded for the winning entry in each of the following categories.

- **Featherweight.** At most 128 words of Logo, not including comments and delimiters.

- **Heavyweight.** At most 1024 words of Logo, not including comments and delimiters.

Winners will be selected by popular vote! (Homework 13)

- Static **images** of the output of your programs

- Tonight at midnight: I'll post your Logo **implementations**!
    - Run them to see these images evolve!

- I will also post a **solution** to the Logo project
    - It runs (almost) all of the contest entries
    - You can use it as a study guide for the final

(Demo)

Tuesday, November 29, 2011

# MapReduce

Tuesday, November 29, 2011

# MapReduce

Tuesday, November 29, 2011

# MapReduce

MapReduce is a *framework* for batch processing of Big Data

Tuesday, November 29, 2011

# MapReduce

MapReduce is a *framework* for batch processing of Big Data

What does that mean?

Tuesday, November 29, 2011

# MapReduce

MapReduce is a *framework* for batch processing of Big Data

What does that mean?

• **Framework:** A system used by programmers to build applications

Tuesday, November 29, 2011

# MapReduce

MapReduce is a *framework* for batch processing of Big Data

What does that mean?

- **Framework:** A system used by programmers to build applications

- **Batch processing:** All the data is available at the outset and results aren't consumed until processing completes

Tuesday, November 29, 2011

# MapReduce

MapReduce is a *framework* for batch processing of Big Data

What does that mean?

- **Framework:** A system used by programmers to build applications

- **Batch processing:** All the data is available at the outset and results aren't consumed until processing completes

- **Big Data:** A buzzword used to describe datasets so large that they reveal facts about the world via statistical analysis

Tuesday, November 29, 2011

# MapReduce

MapReduce is a *framework* for batch processing of Big Data

What does that mean?

- **Framework:** A system used by programmers to build applications

- **Batch processing:** All the data is available at the outset and results aren't consumed until processing completes

- **Big Data:** A buzzword used to describe datasets so large that they reveal facts about the world via statistical analysis

(Demo)

Tuesday, November 29, 2011

# MapReduce

MapReduce is a *framework* for batch processing of Big Data

What does that mean?

- **Framework:** A system used by programmers to build applications

- **Batch processing:** All the data is available at the outset and results aren't consumed until processing completes

- **Big Data:** A buzzword used to describe datasets so large that they reveal facts about the world via statistical analysis

(Demo)

The big ideas that underly MapReduce:

# MapReduce

MapReduce is a *framework* for batch processing of Big Data

What does that mean?

- **Framework:** A system used by programmers to build applications

- **Batch processing:** All the data is available at the outset and results aren't consumed until processing completes

- **Big Data:** A buzzword used to describe datasets so large that they reveal facts about the world via statistical analysis

(Demo)

The big ideas that underly MapReduce:

- Datasets are too big to be stored or analyzed on one machine

Tuesday, November 29, 2011

# MapReduce

MapReduce is a *framework* for batch processing of Big Data

What does that mean?

- **Framework:** A system used by programmers to build applications

- **Batch processing:** All the data is available at the outset and results aren't consumed until processing completes

- **Big Data:** A buzzword used to describe datasets so large that they reveal facts about the world via statistical analysis

(Demo)

The big ideas that underly MapReduce:

- Datasets are too big to be stored or analyzed on one machine

- When using multiple machines, systems issues abound

Tuesday, November 29, 2011

# MapReduce

MapReduce is a *framework* for batch processing of Big Data

What does that mean?

- **Framework:** A system used by programmers to build applications

- **Batch processing:** All the data is available at the outset and results aren't consumed until processing completes

- **Big Data:** A buzzword used to describe datasets so large that they reveal facts about the world via statistical analysis

(Demo)

The big ideas that underly MapReduce:

- Datasets are too big to be stored or analyzed on one machine

- When using multiple machines, systems issues abound

- Pure functions enable an abstraction barrier between data processing logic and distributed system administration

Tuesday, November 29, 2011

# Systems

Tuesday, November 29, 2011

# Systems

Systems research enables the development of applications by defining and implementing abstractions:

# Systems

Systems research enables the development of applications by defining and implementing abstractions:

- **Operating systems** provide a stable, consistent interface to unreliable, inconsistent hardware

Tuesday, November 29, 2011

# Systems

Systems research enables the development of applications by defining and implementing abstractions:

- **Operating systems** provide a stable, consistent interface to unreliable, inconsistent hardware

- **Networks** provide a simple, robust data transfer interface to constantly evolving communications infrastructure

Tuesday, November 29, 2011

# Systems

Systems research enables the development of applications by defining and implementing abstractions:

- **Operating systems** provide a stable, consistent interface to unreliable, inconsistent hardware

- **Networks** provide a simple, robust data transfer interface to constantly evolving communications infrastructure

- **Databases** provide a declarative interface to software that stores and retrieves information efficiently

# Systems

Systems research enables the development of applications by defining and implementing abstractions:

- **Operating systems** provide a stable, consistent interface to unreliable, inconsistent hardware

- **Networks** provide a simple, robust data transfer interface to constantly evolving communications infrastructure

- **Databases** provide a declarative interface to software that stores and retrieves information efficiently

- **Distributed systems** provide a single-entity-level interface to a cluster of multiple machines

# Systems

Systems research enables the development of applications by defining and implementing abstractions:

- **Operating systems** provide a stable, consistent interface to unreliable, inconsistent hardware

- **Networks** provide a simple, robust data transfer interface to constantly evolving communications infrastructure

- **Databases** provide a declarative interface to software that stores and retrieves information efficiently

- **Distributed systems** provide a single-entity-level interface to a cluster of multiple machines

A unifying property of effective systems:

# Systems

Systems research enables the development of applications by defining and implementing abstractions:

- **Operating systems** provide a stable, consistent interface to unreliable, inconsistent hardware

- **Networks** provide a simple, robust data transfer interface to constantly evolving communications infrastructure

- **Databases** provide a declarative interface to software that stores and retrieves information efficiently

- **Distributed systems** provide a single-entity-level interface to a cluster of multiple machines

A unifying property of effective systems:

Hide *complexity*, but retain *flexibility*

Tuesday, November 29, 2011

# The Unix Operating System

Tuesday, November 29, 2011

# The Unix Operating System

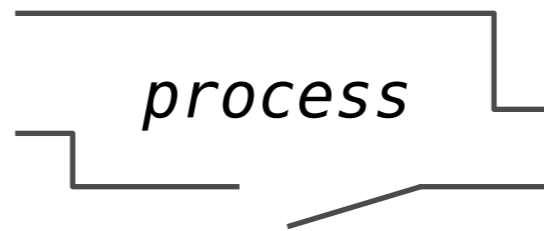Essential features of the Unix operating system (and variants)

# The Unix Operating System

Essential features of the Unix operating system (and variants)

- **Portability:** The same operating system on different hardware

# The Unix Operating System
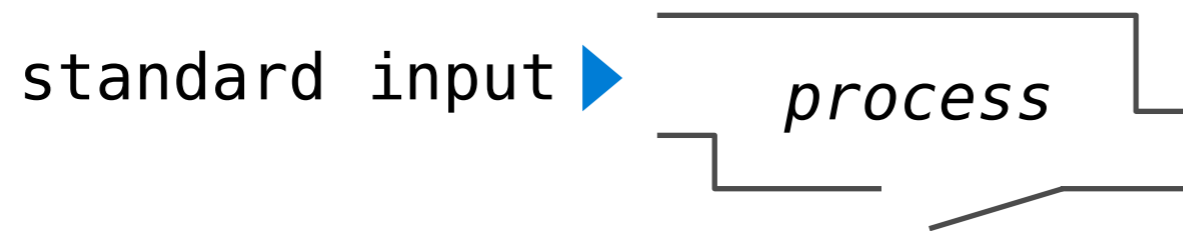
Essential features of the Unix operating system (and variants)

- **Portability:** The same operating system on different hardware

- **Multi-Tasking:** Many processes run concurrently on a machine

# The Unix Operating System

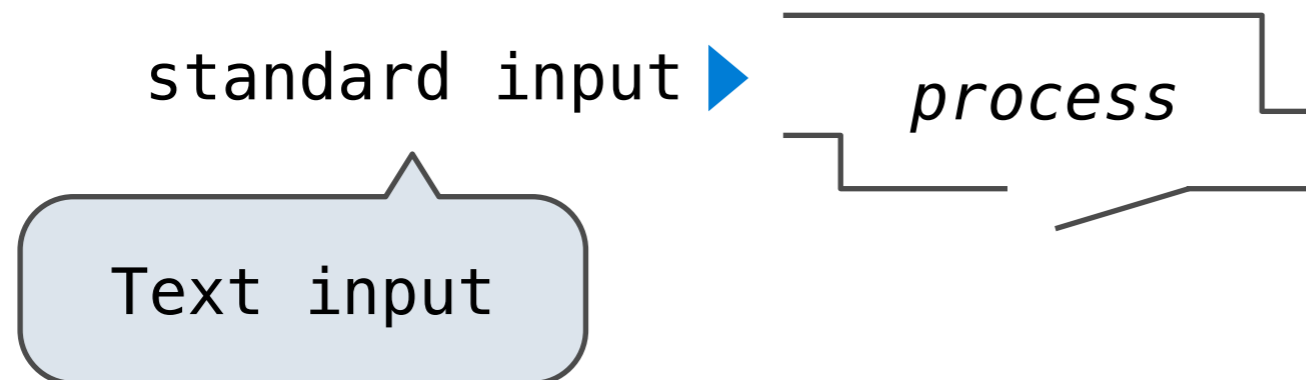Essential features of the Unix operating system (and variants)

- **Portability:** The same operating system on different hardware

- **Multi-Tasking:** Many processes run concurrently on a machine

- **Plain Text:** Data is stored and shared in text format

# The Unix Operating System

Essential features of the Unix operating system (and variants)

- **Portability:** The same operating system on different hardware

- **Multi-Tasking:** Many processes run concurrently on a machine

- **Plain Text:** Data is stored and shared in text format

- **Modularity:** Small tools are composed flexibly via pipes

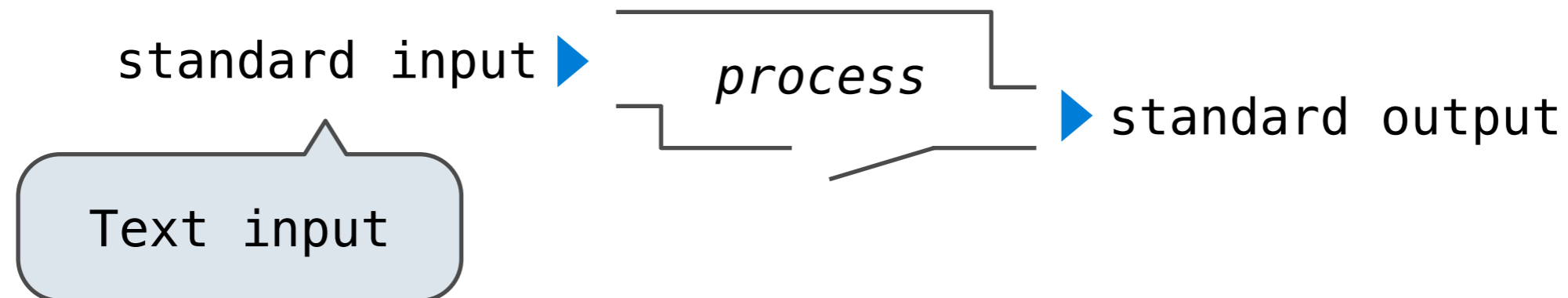Essential features of the Unix operating system (and variants)

- **Portability:** The same operating system on different hardware

- **Multi-Tasking:** Many processes run concurrently on a machine

- **Plain Text:** Data is stored and shared in text format

- **Modularity:** Small tools are composed flexibly via pipes

*process*

# The Unix Operating System

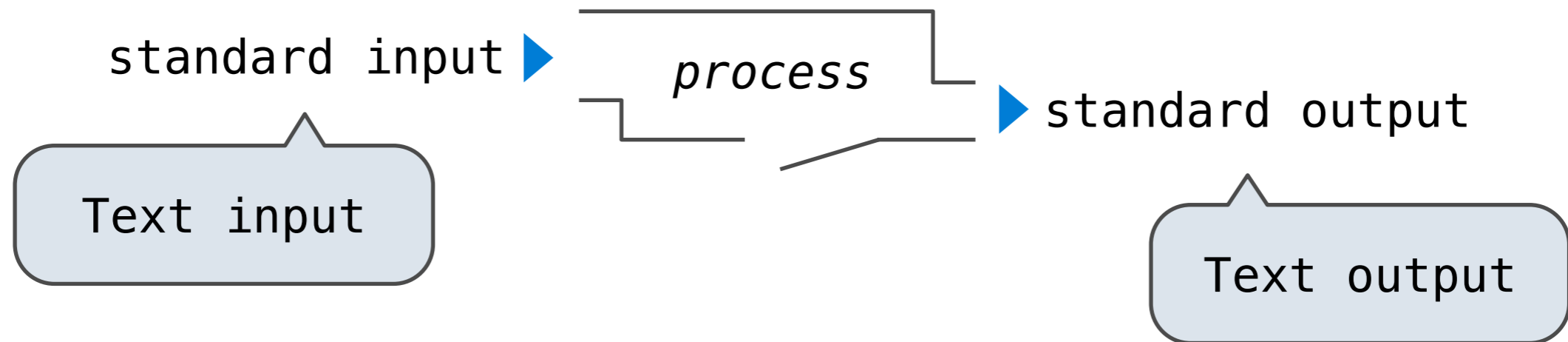Essential features of the Unix operating system (and variants)

- **Portability:** The same operating system on different hardware

- **Multi-Tasking:** Many processes run concurrently on a machine

- **Plain Text:** Data is stored and shared in text format

- **Modularity:** Small tools are composed flexibly via pipes

standard input ▶ *process*

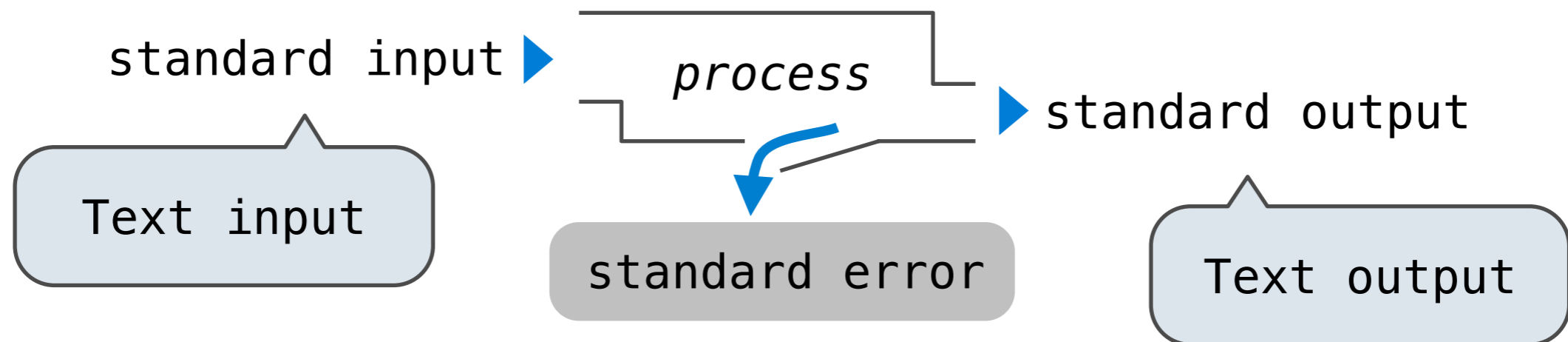Essential features of the Unix operating system (and variants)

- **Portability:** The same operating system on different hardware
- **Multi-Tasking:** Many processes run concurrently on a machine
- **Plain Text:** Data is stored and shared in text format
- **Modularity:** Small tools are composed flexibly via pipes

standard input ▶ *process*

Text input

# The Unix Operating System

Essential features of the Unix operating system (and variants)

- **Portability:** The same operating system on different hardware

- **Multi-Tasking:** Many processes run concurrently on a machine

- **Plain Text:** Data is stored and shared in text format

- **Modularity:** Small tools are composed flexibly via pipes

standard input ▶ *process* ▶ standard output

Text input

# The Unix Operating System

Essential features of the Unix operating system (and variants)
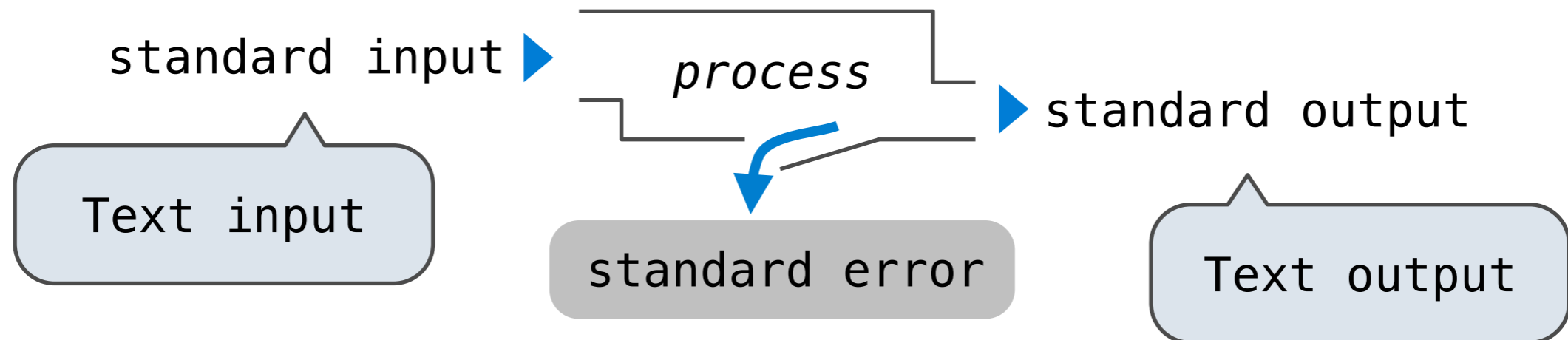
- **Portability:** The same operating system on different hardware
- **Multi-Tasking:** Many processes run concurrently on a machine
- **Plain Text:** Data is stored and shared in text format
- **Modularity:** Small tools are composed flexibly via pipes

standard input ▶ *process* ▶ standard output

Text input

Text output

# The Unix Operating System

Essential features of the Unix operating system (and variants)

- **Portability:** The same operating system on different hardware
- **Multi-Tasking:** Many processes run concurrently on a machine
- **Plain Text:** Data is stored and shared in text format
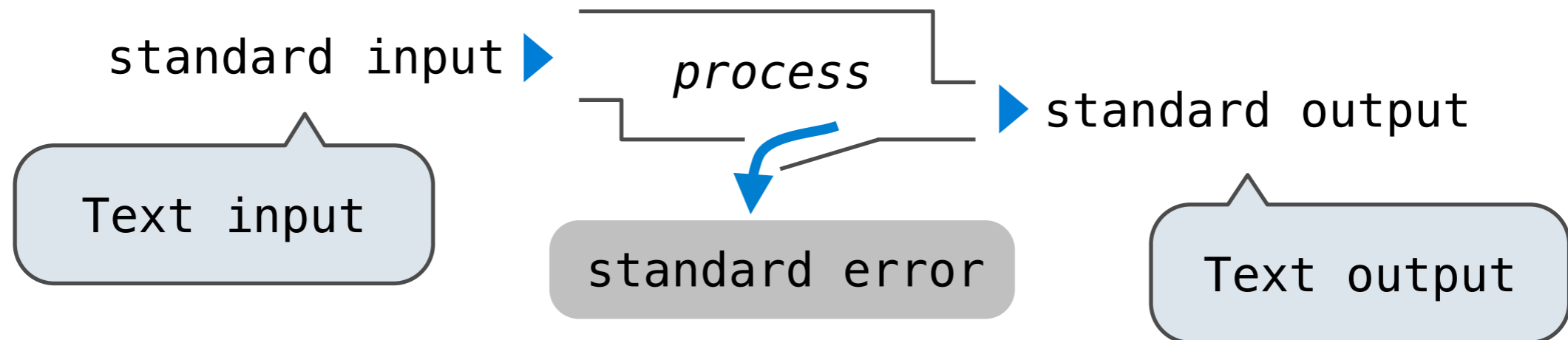- **Modularity:** Small tools are composed flexibly via pipes

standard input ▶ *process* ▶ standard output

Text input

standard error

Text output

# The Unix Operating System

Essential features of the Unix operating system (and variants)

- **Portability:** The same operating system on different hardware
- **Multi-Tasking:** Many processes run concurrently on a machine
- **Plain Text:** Data is stored and shared in text format
- **Modularity:** Small tools are composed flexibly via pipes

standard input ▶ *process* ▶ standard output

Text input

standard error

Text output

The ***standard streams*** in a Unix-like operating system

are conceptually similar to Python iterators

# The Unix Operating System

Essential features of the Unix operating system (and variants)

- **Portability:** The same operating system on different hardware
- **Multi-Tasking:** Many processes run concurrently on a machine
- **Plain Text:** Data is stored and shared in text format
- **Modularity:** Small tools are composed flexibly via pipes

standard input ▶ *process* ▶ standard output

Text input

standard error

Text output

The ***standard streams*** in a Unix-like operating system

are conceptually similar to Python iterators

(Demo)

# Python Programs in a Unix Environment

# Python Programs in a Unix Environment

The built-in `input` function reads a line from *standard input*.

# Python Programs in a Unix Environment

The built-in `input` function reads a line from *standard input*.

The built-in `print` function writes a line to *standard output*.

The built-in `input` function reads a line from *standard input*.

The built-in `print` function writes a line to *standard output*.

(Demo)

# Python Programs in a Unix Environment

The built-in `input` function reads a line from *standard input*.

The built-in `print` function writes a line to *standard output*.

(Demo)

The values `sys.stdin` and `sys.stdout` also provide access to the Unix *standard streams* as "files."

# Python Programs in a Unix Environment

The built-in `input` function reads a line from *standard input*.

The built-in `print` function writes a line to *standard output*.

(Demo)

The values `sys.stdin` and `sys.stdout` also provide access to the Unix *standard streams* as "files."

A Python "file" is an interface that supports iteration, read, and write methods.

# Python Programs in a Unix Environment

The built-in `input` function reads a line from *standard input*.

The built-in `print` function writes a line to *standard output*.

(Demo)

The values `sys.stdin` and `sys.stdout` also provide access to the Unix *standard streams* as "files."

A Python "file" is an interface that supports iteration, read, and write methods.

Using these "files" takes advantage of the operating system *standard stream* abstraction.

Tuesday, November 29, 2011

# Python Programs in a Unix Environment

The built-in `input` function reads a line from *standard input*.

The built-in `print` function writes a line to *standard output*.

(Demo)

The values `sys.stdin` and `sys.stdout` also provide access to the Unix *standard streams* as "files."

A Python "file" is an interface that supports iteration, read, and write methods.

Using these "files" takes advantage of the operating system *standard stream* abstraction.

(Demo)

Tuesday, November 29, 2011

# MapReduce Evaluation Model

Tuesday, November 29, 2011

# MapReduce Evaluation Model

**Map phase:** Apply a *mapper* function to inputs, emitting a set of intermediate key-value pairs

# MapReduce Evaluation Model

**Map phase:** Apply a *mapper* function to inputs, emitting a set of intermediate key-value pairs

- The *mapper* takes an iterator over inputs, such as text lines.

# MapReduce Evaluation Model

**Map phase:** Apply a *mapper* function to inputs, emitting a set of intermediate key-value pairs

- The *mapper* takes an iterator over inputs, such as text lines.

- The *mapper* yields 0 or more key-value pairs per input.

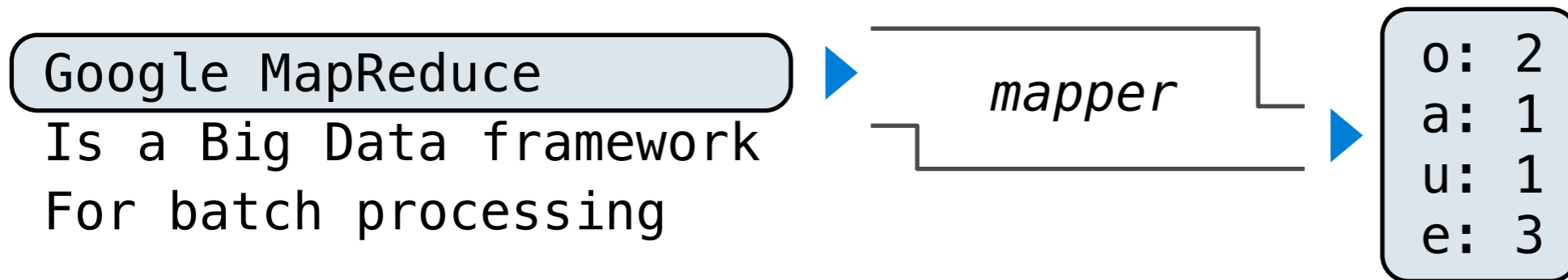# MapReduce Evaluation Model

**Map phase:** Apply a *mapper* function to inputs, emitting a set of intermediate key-value pairs

- The *mapper* takes an iterator over inputs, such as text lines.

- The *mapper* yields 0 or more key-value pairs per input.

```
Google MapReduce
Is a Big Data framework
For batch processing
```

# MapReduce Evaluation Model

**Map phase:** Apply a *mapper* function to inputs, emitting a set of intermediate key-value pairs

- The *mapper* takes an iterator over inputs, such as text lines.

- The *mapper* yields 0 or more key-value pairs per input.

Google MapReduce
Is a Big Data framework
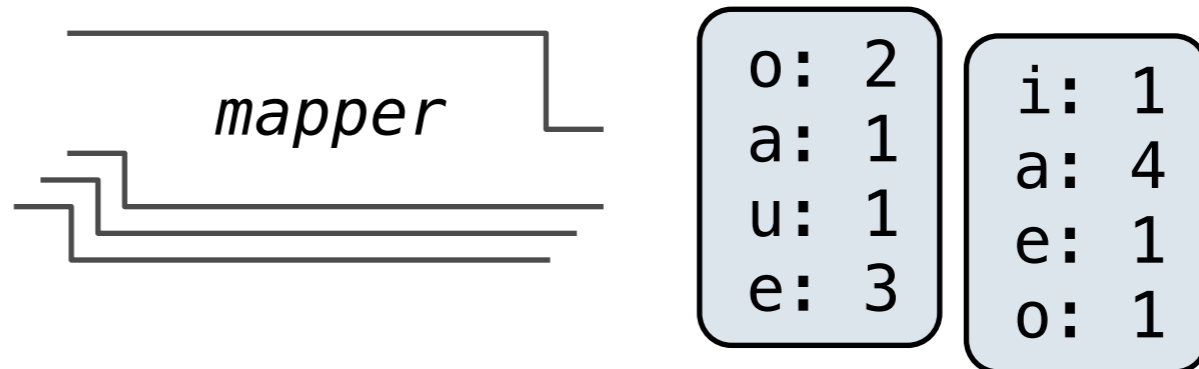For batch processing

*mapper*

# MapReduce Evaluation Model

**Map phase**: Apply a *mapper* function to inputs, emitting a set of intermediate key-value pairs

- The *mapper* takes an iterator over inputs, such as text lines.

- The *mapper* yields 0 or more key-value pairs per input.

Google MapReduce
Is a Big Data framework
For batch processing

*mapper*

o: 2
a: 1
u: 1
e: 3

# MapReduce Evaluation Model

**Map phase:** Apply a *mapper* function to inputs, emitting a set of intermediate key-value pairs

- The *mapper* takes an iterator over inputs, such as text lines.

- The *mapper* yields 0 or more key-value pairs per input.

```
Google MapReduce
Is a Big Data framework
For batch processing
```
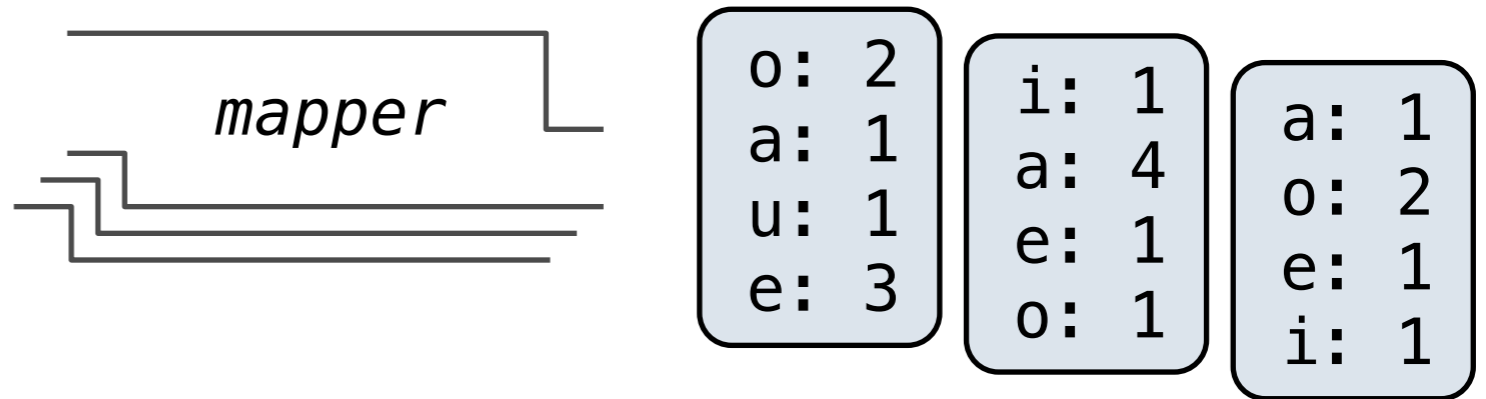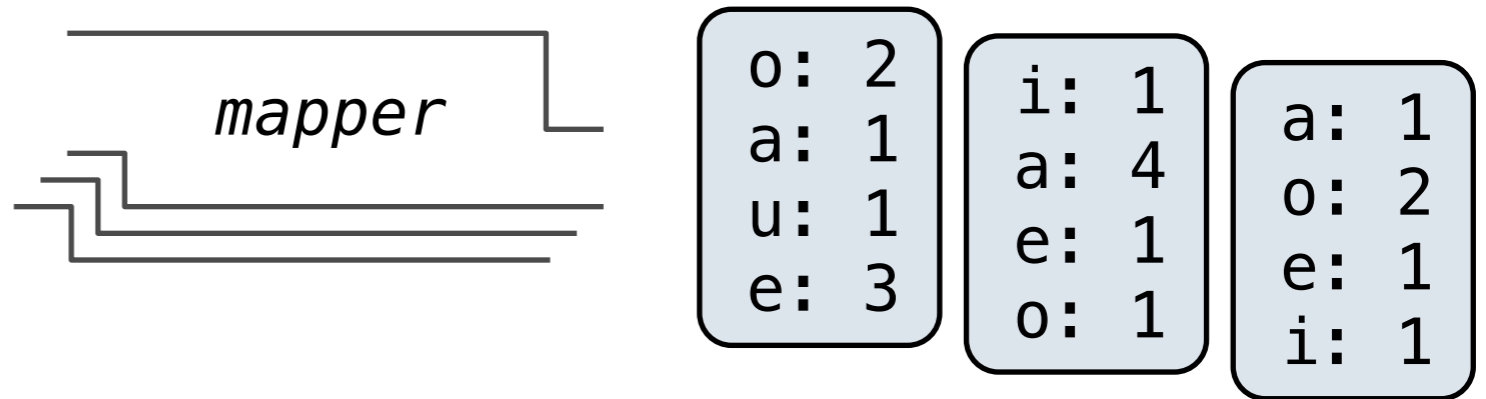
*mapper*

```
o: 2
a: 1
u: 1
e: 3
```

# MapReduce Evaluation Model

**Map phase**: Apply a *mapper* function to inputs, emitting a set of intermediate key-value pairs

- The *mapper* takes an iterator over inputs, such as text lines.

- The *mapper* yields 0 or more key-value pairs per input.

```
Google MapReduce
Is a Big Data framework                mapper
For batch processing
```

```
o: 2        i: 1
a: 1        a: 4
u: 1        e: 1
e: 3        o: 1
```

# MapReduce Evaluation Model

**Map phase**: Apply a *mapper* function to inputs, emitting a set of intermediate key-value pairs

- The *mapper* takes an iterator over inputs, such as text lines.

- The *mapper* yields 0 or more key-value pairs per input.

```
Google MapReduce
Is a Big Data framework
For batch processing
```
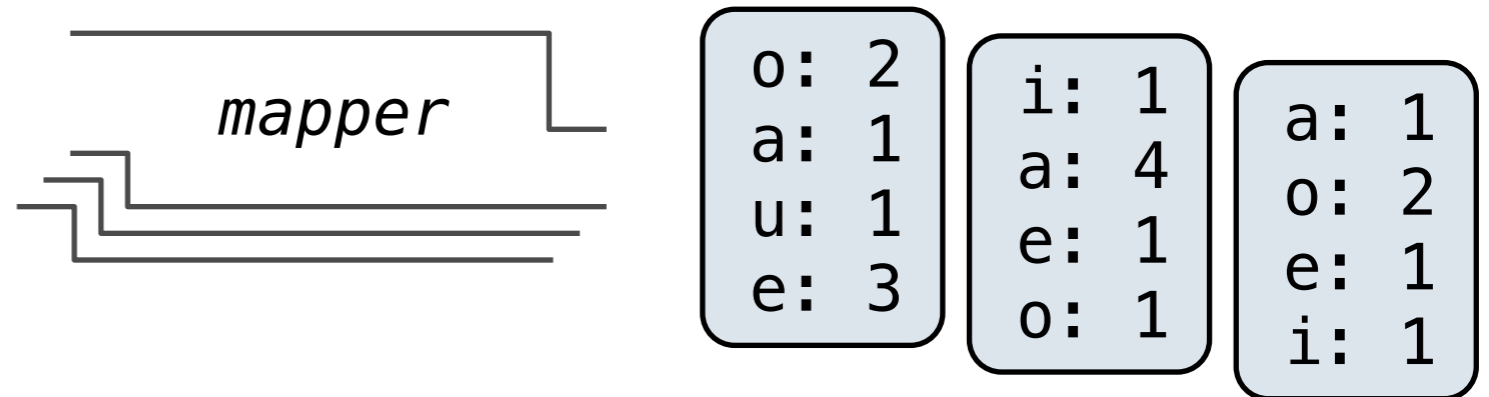
*mapper*

```
o: 2
a: 1
u: 1
e: 3
```

```
i: 1
a: 4
e: 1
o: 1
```

```
a: 1
o: 2
e: 1
i: 1
```

**Map phase:** Apply a *mapper* function to inputs, emitting a set of intermediate key-value pairs

- The *mapper* takes an iterator over inputs, such as text lines.

- The *mapper* yields 0 or more key-value pairs per input.

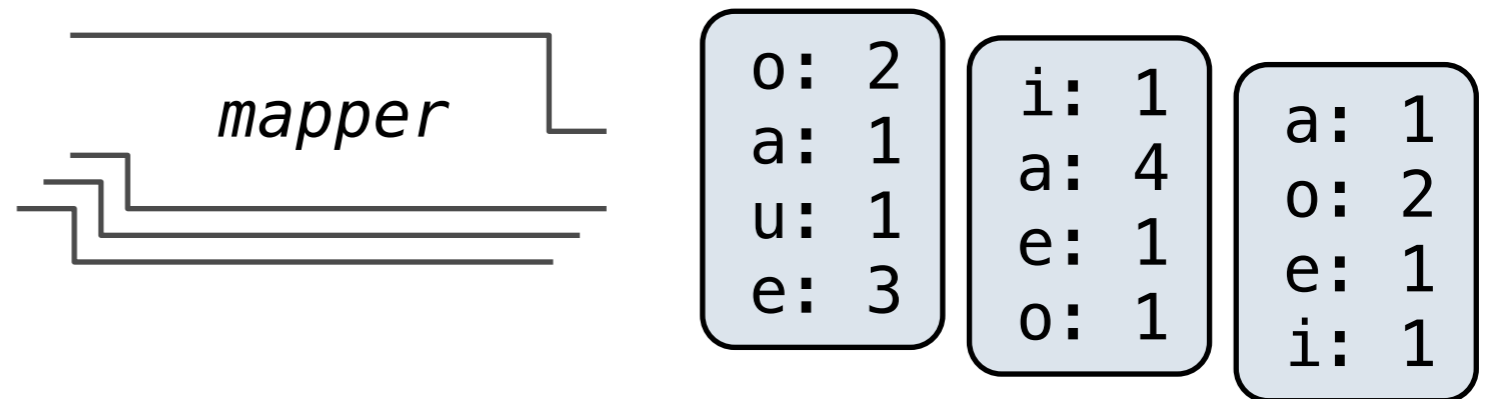Google MapReduce
Is a Big Data framework
For batch processing

*mapper*

| o: 2 |
| a: 1 |
| u: 1 |
| e: 3 |

| i: 1 |
| a: 4 |
| e: 1 |
| o: 1 |

| a: 1 |
| o: 2 |
| e: 1 |
| i: 1 |

Tuesday, November 29, 2011

# MapReduce Evaluation Model

**Map phase:** Apply a *mapper* function to inputs, emitting a set of intermediate key-value pairs

- The *mapper* takes an iterator over inputs, such as text lines.

- The *mapper* yields 0 or more key-value pairs per input.

Google MapReduce
Is a Big Data framework
For batch processing

*mapper*

```
o: 2
a: 1
u: 1
e: 3
```

```
i: 1
a: 4
e: 1
o: 1
```

```
a: 1
o: 2
e: 1
i: 1
```

**Reduce phase:** For each intermediate key, apply a *reducer* function to accumulate all values associated with that key

# MapReduce Evaluation Model

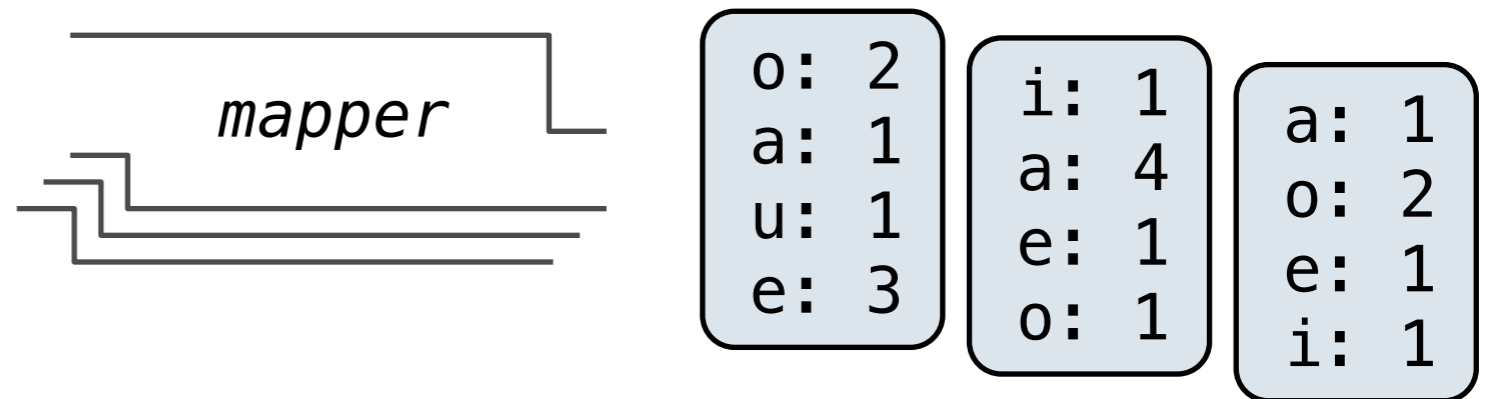**Map phase**: Apply a *mapper* function to inputs, emitting a set of intermediate key-value pairs

- The *mapper* takes an iterator over inputs, such as text lines.

- The *mapper* yields 0 or more key-value pairs per input.

Google MapReduce
Is a Big Data framework
For batch processing

*mapper*

| o: 2 |
| a: 1 |
| u: 1 |
| e: 3 |

| i: 1 |
| a: 4 |
| e: 1 |
| o: 1 |

| a: 1 |
| o: 2 |
| e: 1 |
| i: 1 |

**Reduce phase**: For each intermediate key, apply a *reducer* function to accumulate all values associated with that key

- The *reducer takes* an iterator over key-value pairs.
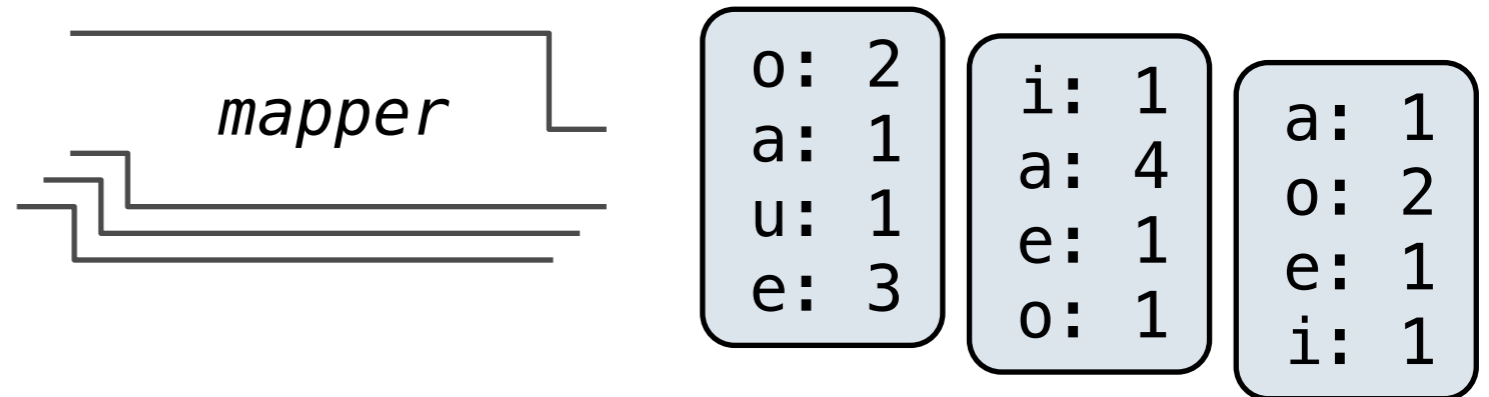
# MapReduce Evaluation Model

**Map phase:** Apply a *mapper* function to inputs, emitting a set of intermediate key-value pairs

• The *mapper* takes an iterator over inputs, such as text lines.

• The *mapper* yields 0 or more key-value pairs per input.

Google MapReduce
Is a Big Data framework
For batch processing

*mapper*

| o: 2 |
| a: 1 |
| u: 1 |
| e: 3 |

| i: 1 |
| a: 4 |
| e: 1 |
| o: 1 |

| a: 1 |
| o: 2 |
| e: 1 |
| i: 1 |

**Reduce phase:** For each intermediate key, apply a *reducer* function to accumulate all values associated with that key

• The *reducer takes* an iterator over key-value pairs.

• All pairs with a given key are consecutive

# MapReduce Evaluation Model

**Map phase:** Apply a *mapper* function to inputs, emitting a set of intermediate key-value pairs

- The *mapper* takes an iterator over inputs, such as text lines.

- The *mapper* yields 0 or more key-value pairs per input.

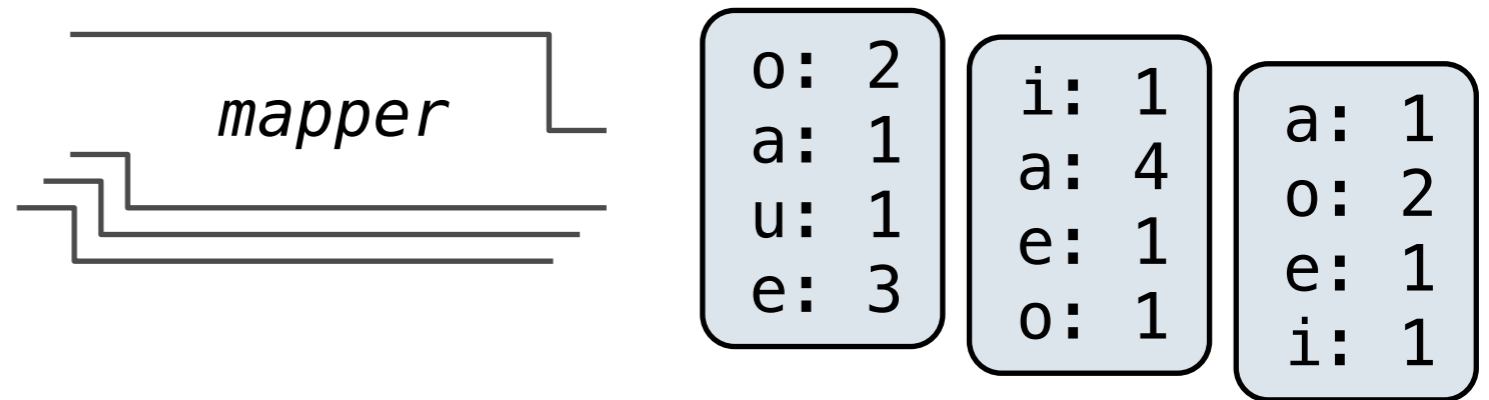Google MapReduce
Is a Big Data framework
For batch processing

*mapper*

| o: 2 |
| a: 1 |
| u: 1 |
| e: 3 |

| i: 1 |
| a: 4 |
| e: 1 |
| o: 1 |

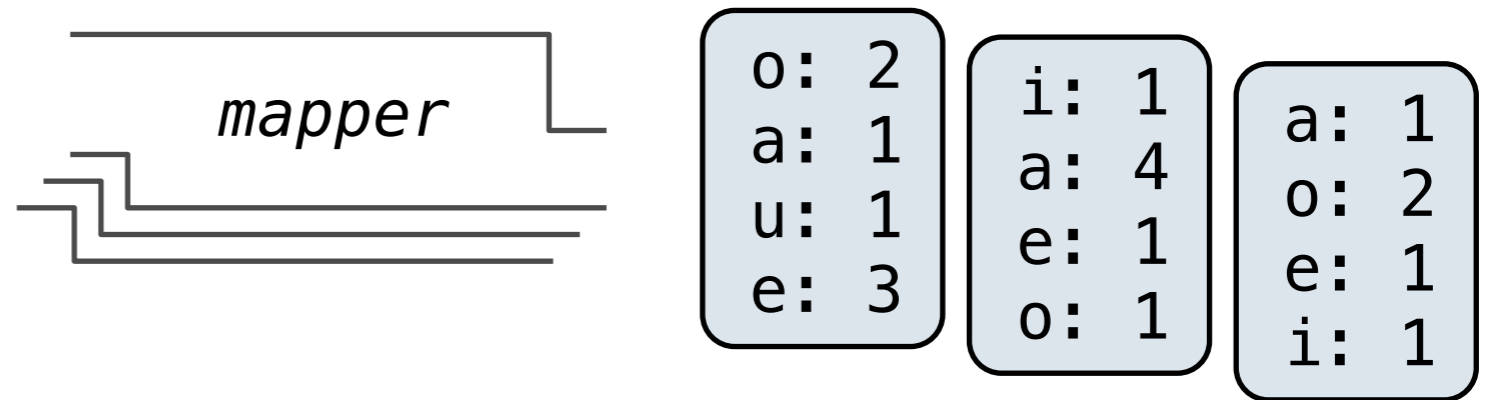| a: 1 |
| o: 2 |
| e: 1 |
| i: 1 |

**Reduce phase:** For each intermediate key, apply a *reducer* function to accumulate all values associated with that key

- The *reducer takes* an iterator over key-value pairs.

- All pairs with a given key are consecutive

- The *reducer* yields 0 or more values for a key, each associated with that intermediate key.

# MapReduce Evaluation Model

```
Google MapReduce
Is a Big Data framework
For batch processing
```

*mapper*

```
o: 2
a: 1
u: 1
e: 3
```

```
i: 1
a: 4
e: 1
o: 1
```

```
a: 1
o: 2
e: 1
i: 1
```

**Reduce phase:** For each intermediate key, apply a *reducer* function to accumulate all values associated with that key

• The *reducer takes* an iterator over key-value pairs.

• All pairs with a given key are consecutive

• The *reducer* yields 0 or more values for a key, each associated with that intermediate key.

# MapReduce Evaluation Model

Google MapReduce
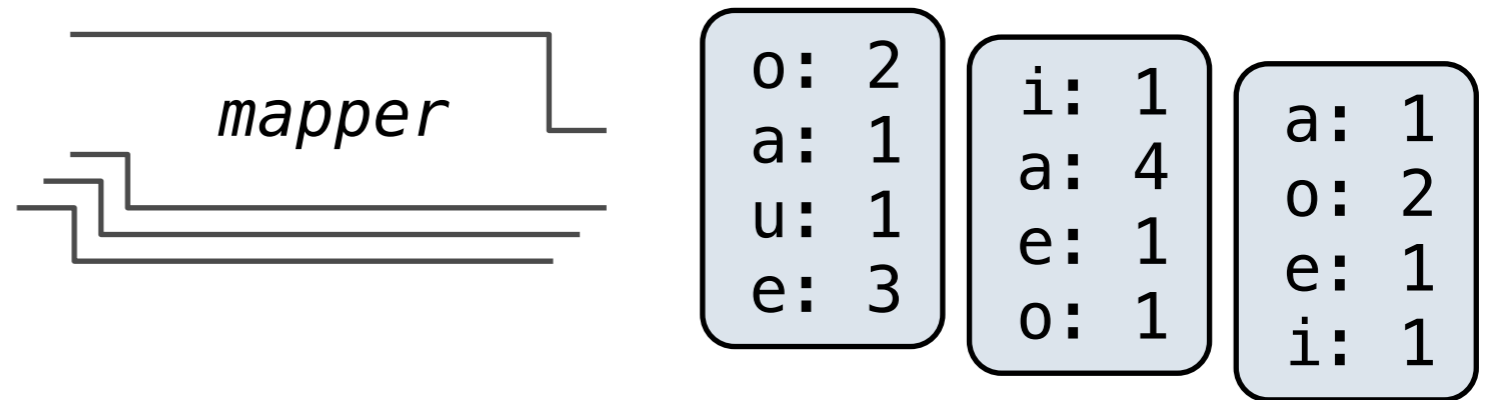Is a Big Data framework
For batch processing

*mapper*

```
o: 2
a: 1
u: 1
e: 3
```

```
i: 1
a: 4
e: 1
o: 1
```

```
a: 1
o: 2
e: 1
i: 1
```

**Reduce phase:** For each intermediate key, apply a *reducer* function to accumulate all values associated with that key

- The *reducer takes* an iterator over key-value pairs.

- All pairs with a given key are consecutive

- The *reducer* yields 0 or more values for a key, each associated with that intermediate key.
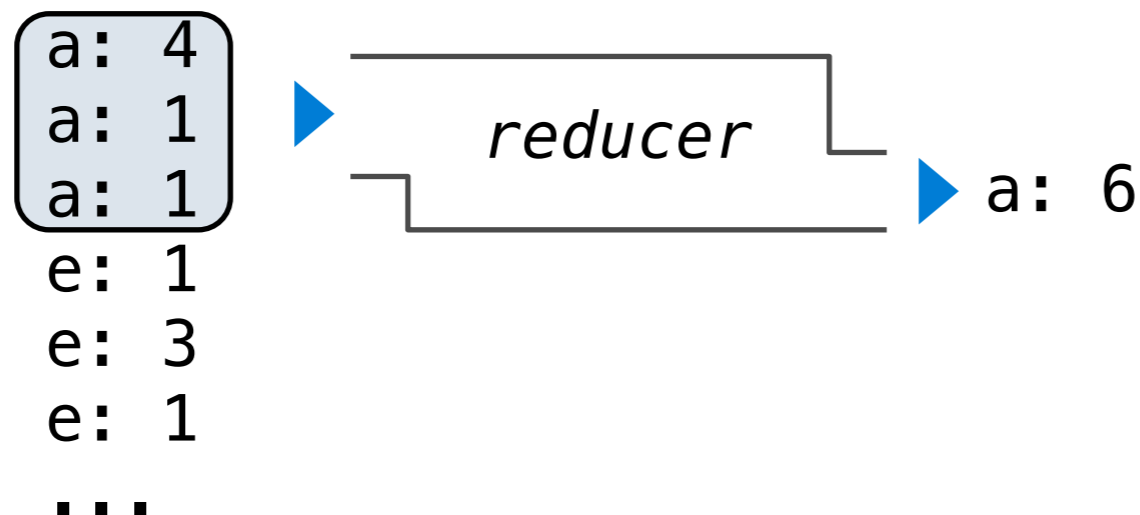
```
a: 4
a: 1
a: 1
e: 1
e: 3
e: 1
...
```

# MapReduce Evaluation Model

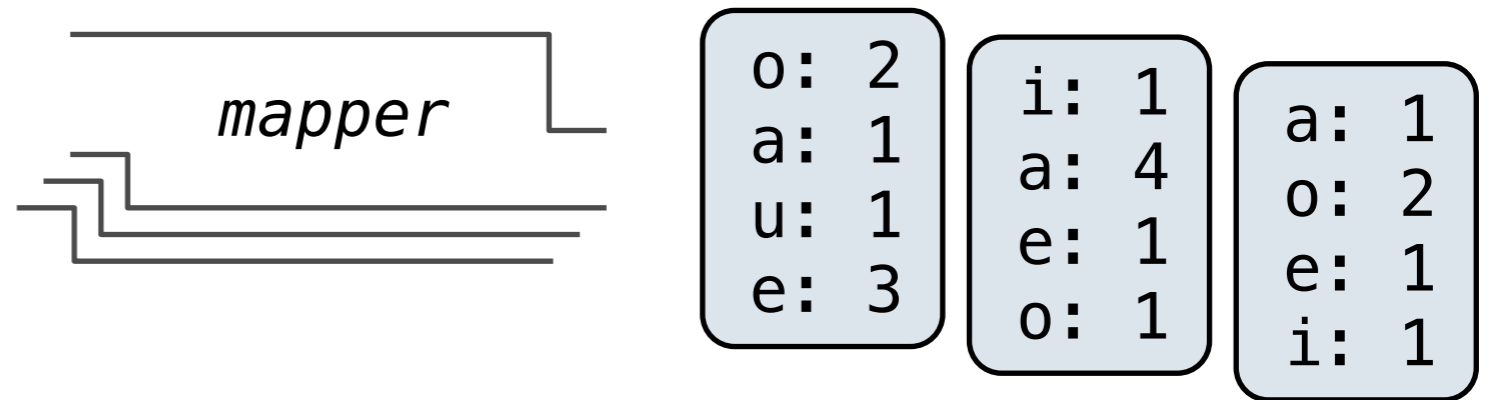Google MapReduce
Is a Big Data framework
For batch processing

*mapper*

| o: 2 |
| a: 1 |
| u: 1 |
| e: 3 |

| i: 1 |
| a: 4 |
| e: 1 |
| o: 1 |

| a: 1 |
| o: 2 |
| e: 1 |
| i: 1 |

**Reduce phase:** For each intermediate key, apply a *reducer* function to accumulate all values associated with that key

- The *reducer takes* an iterator over key-value pairs.

- All pairs with a given key are consecutive

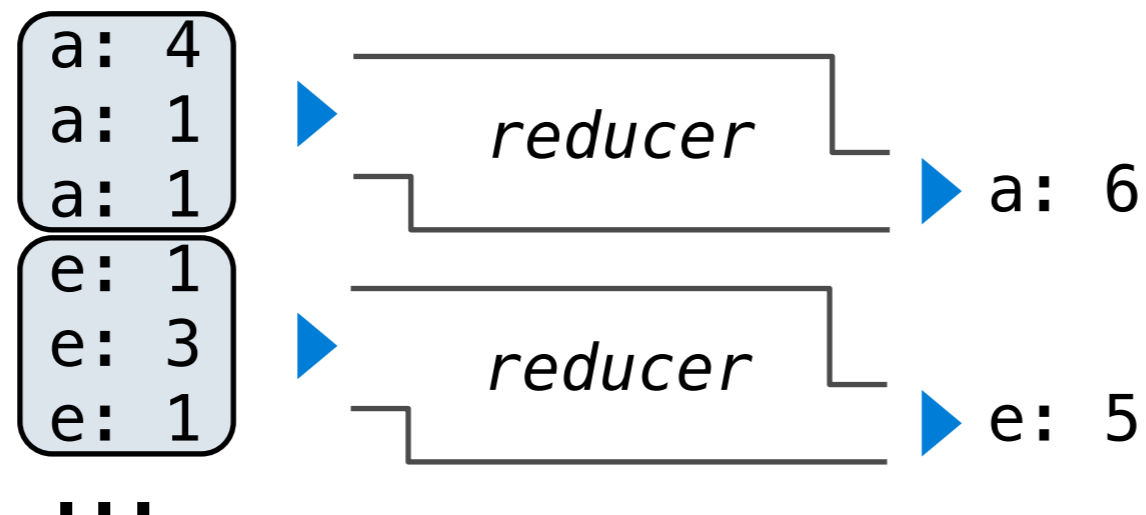- The *reducer* yields 0 or more values for a key, each associated with that intermediate key.

| a: 4 |
| a: 1 |
| a: 1 |
| e: 1 |
| e: 3 |
| e: 1 |
| ... |

▶ *reducer* ▶ a: 6

Tuesday, November 29, 2011

# MapReduce Evaluation Model

Google MapReduce
Is a Big Data framework
For batch processing

*mapper*

```
o: 2
a: 1
u: 1
e: 3
```

```
i: 1
a: 4
e: 1
o: 1
```
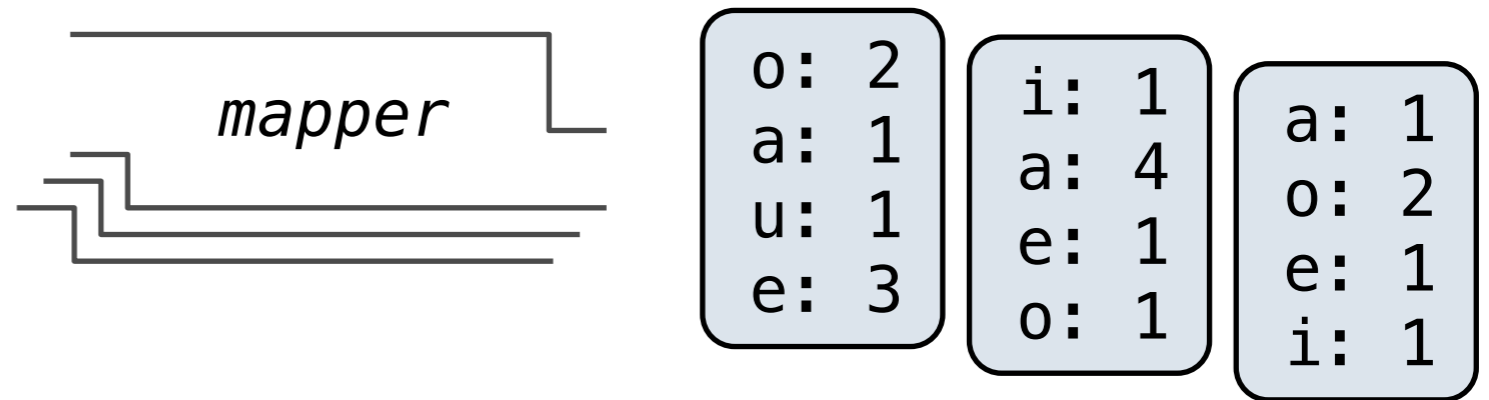
```
a: 1
o: 2
e: 1
i: 1
```

**Reduce phase:** For each intermediate key, apply a *reducer* function to accumulate all values associated with that key

- The *reducer takes* an iterator over key-value pairs.

- All pairs with a given key are consecutive

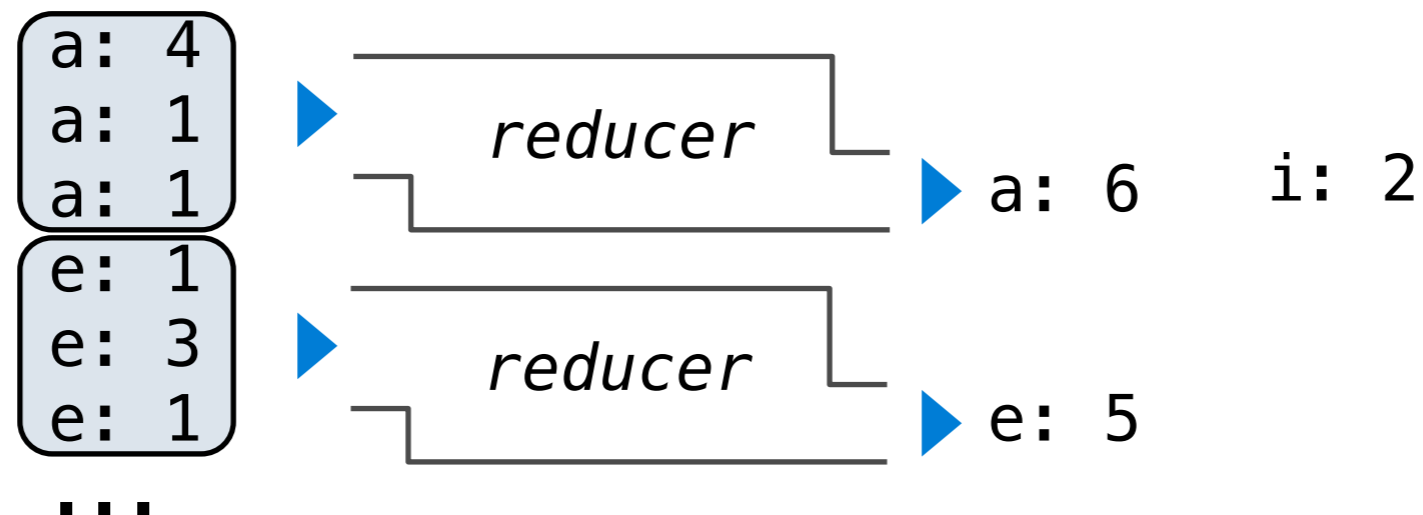- The *reducer* yields 0 or more values for a key, each associated with that intermediate key.

```
a: 4
a: 1
a: 1
```

▶ *reducer* ▶ a: 6

```
e: 1
e: 3
e: 1
```

▶ *reducer* ▶ e: 5

...

# MapReduce Evaluation Model

Google MapReduce
Is a Big Data framework
For batch processing

*mapper*

```
o: 2        i: 1        a: 1
a: 1        a: 4        o: 2
u: 1        e: 1        e: 1
e: 3        o: 1        i: 1
```
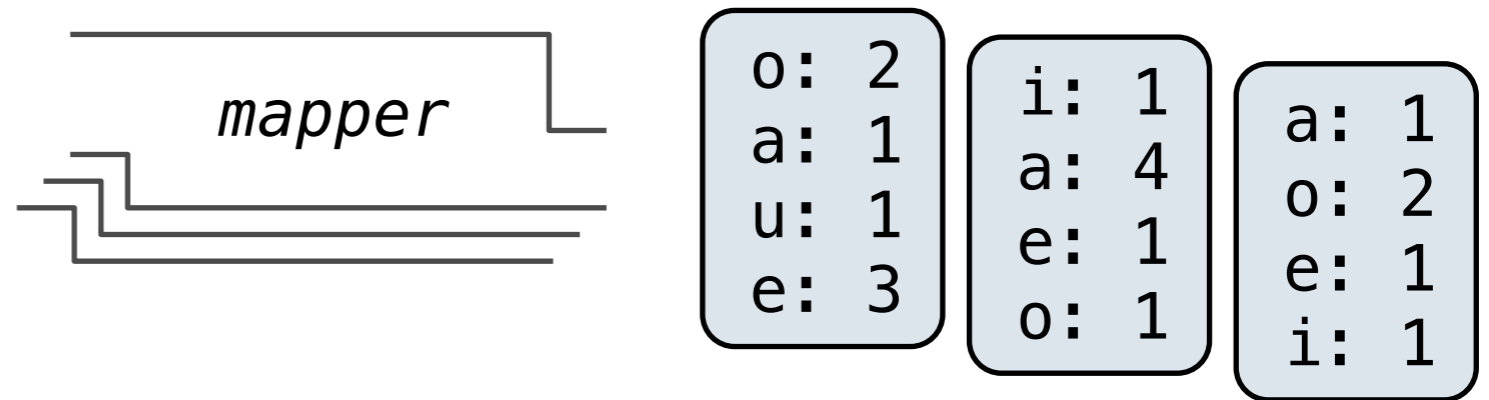
**Reduce phase:** For each intermediate key, apply a *reducer* function to accumulate all values associated with that key

- The *reducer takes* an iterator over key-value pairs.

- All pairs with a given key are consecutive

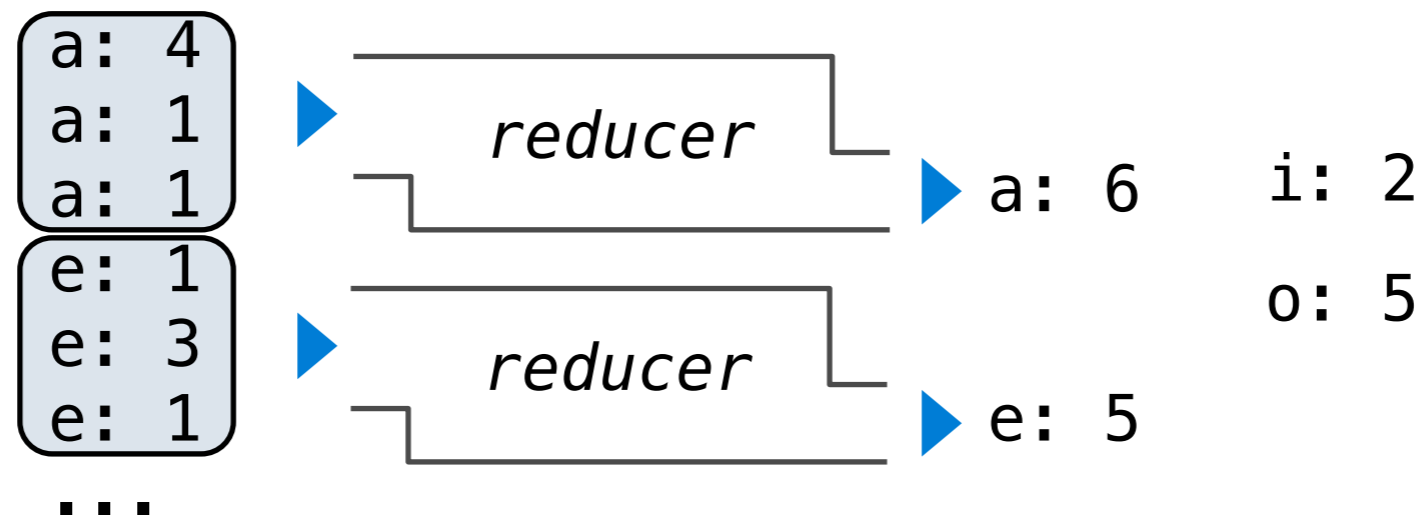- The *reducer* yields 0 or more values for a key, each associated with that intermediate key.

```
a: 4
a: 1    ▶   reducer
a: 1                    ▶  a: 6      i: 2
e: 1
e: 3    ▶   reducer
e: 1                    ▶  e: 5
...
```

# MapReduce Evaluation Model

Google MapReduce
Is a Big Data framework
For batch processing

*mapper*

```
o: 2
a: 1
u: 1
e: 3
```

```
i: 1
a: 4
e: 1
o: 1
```
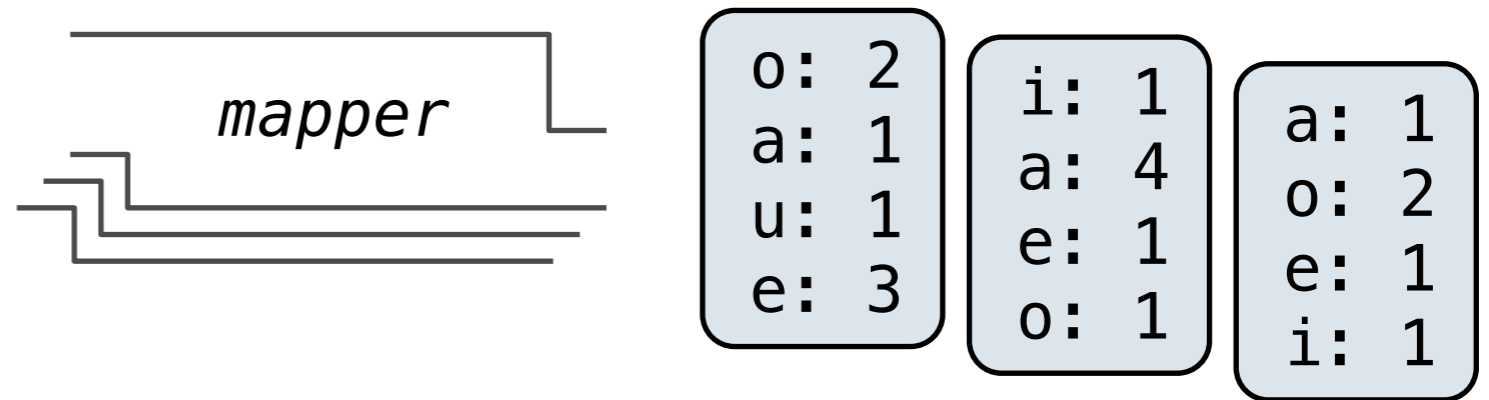
```
a: 1
o: 2
e: 1
i: 1
```

**Reduce phase:** For each intermediate key, apply a *reducer* function to accumulate all values associated with that key

- The *reducer takes* an iterator over key-value pairs.

- All pairs with a given key are consecutive

- The *reducer* yields 0 or more values for a key, each associated with that intermediate key.

```
a: 4
a: 1
a: 1
```

*reducer*

a: 6

```
e: 1
e: 3
e: 1
```

*reducer*

e: 5

i: 2

o: 5

...

8

Tuesday, November 29, 2011

# MapReduce Evaluation Model

Google MapReduce
Is a Big Data framework
For batch processing

*mapper*

```
o: 2     i: 1     a: 1
a: 1     a: 4     o: 2
u: 1     e: 1     e: 1
e: 3     o: 1     i: 1
```

**Reduce phase:** For each intermediate key, apply a *reducer* function to accumulate all values associated with that key

- The *reducer takes* an iterator over key-value pairs.

- All pairs with a given key are consecutive

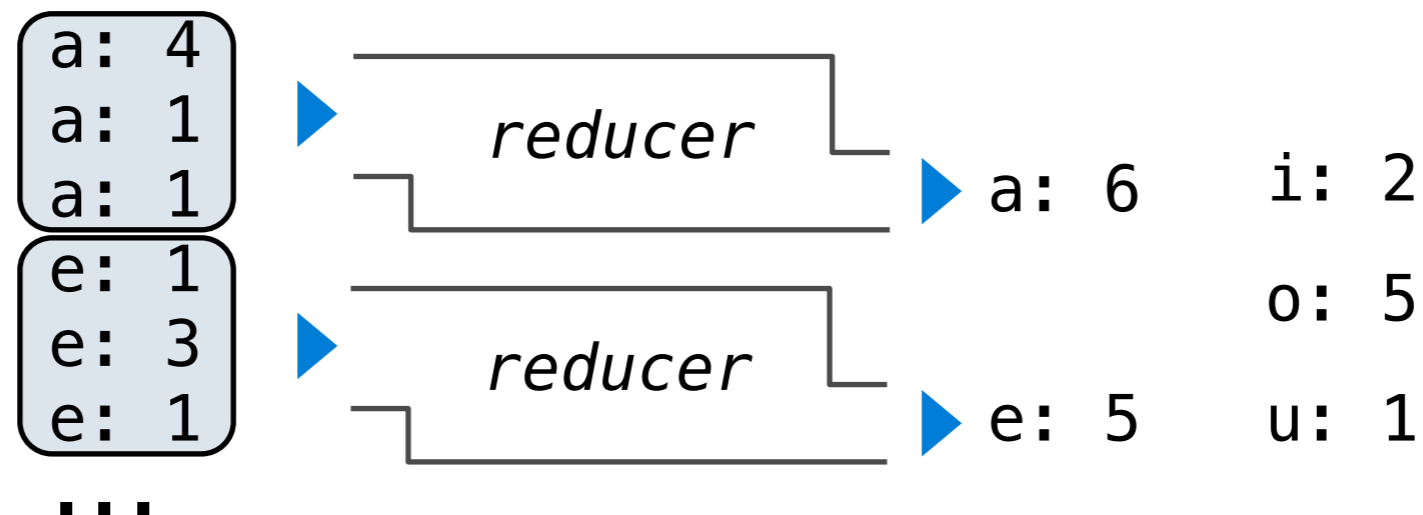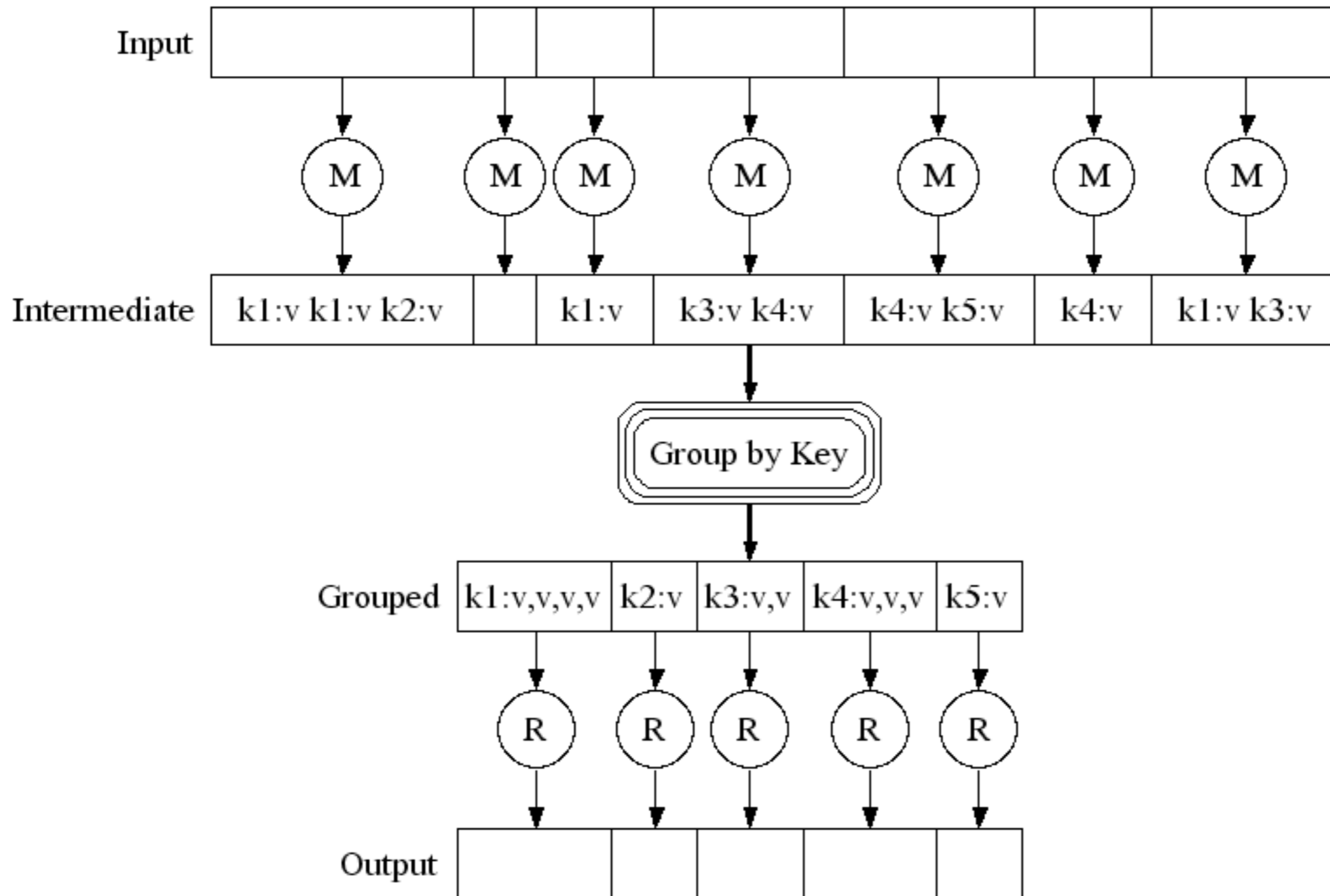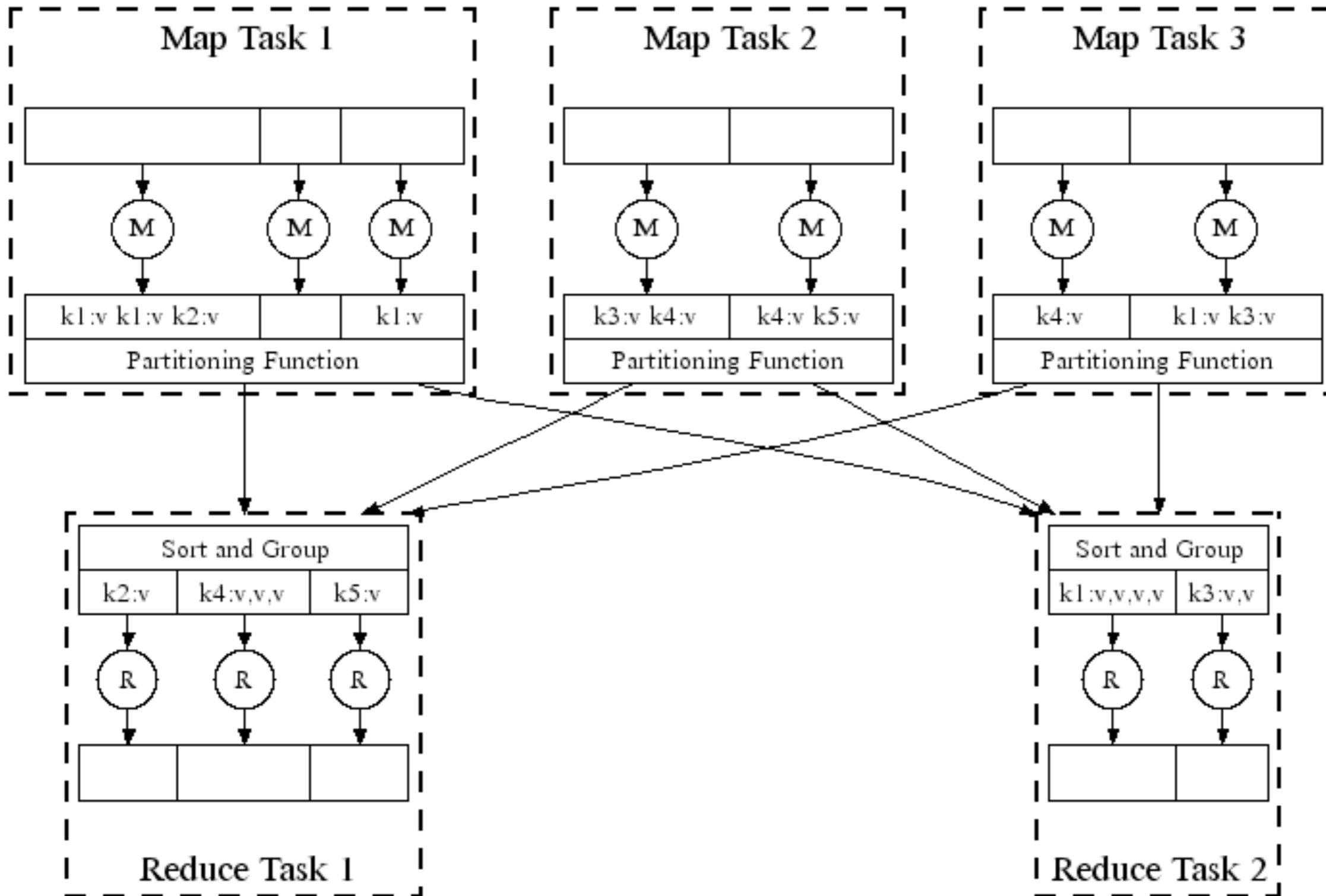- The *reducer* yields 0 or more values for a key, each associated with that intermediate key.

```
a: 4
a: 1     ▶   reducer          ▶  a: 6      i: 2
a: 1
e: 1                                        o: 5
e: 3     ▶   reducer
e: 1                          ▶  e: 5      u: 1
...
```

# Above-the-Line: Execution model

Tuesday, November 29, 2011

# Below-the-Line: Parallel Execution

http://research.google.com/archive/mapreduce-osdi04-slides/index-auto-0008.html

Tuesday, November 29, 2011

# Below-the-Line: Parallel Execution



A "task" is a Unix process running on a machine

Tuesday, November 29, 2011

# Below-the-Line: Parallel Execution



A "task" is a Unix process running on a machine

# Python Examples of a MapReduce Application

# Python Examples of a MapReduce Application

The *mapper* and *reducer* are both self-contained Python programs

# Python Examples of a MapReduce Application

The *mapper* and *reducer* are both self-contained Python programs

- Read from *standard input* and write to *standard output*!

The *mapper* and *reducer* are both self-contained Python programs

- Read from *standard input* and write to *standard output*!

**Mapper**

# Python Examples of a MapReduce Application

The *mapper* and *reducer* are both self-contained Python programs

- Read from *standard input* and write to *standard output*!

**Mapper**

```python
def emit_vowels(line):
    for vowel in 'aeiou':
        count = line.count(vowel)
        if count > 0:
            emit(vowel, count)
```

# Python Examples of a MapReduce Application

The *mapper* and *reducer* are both self-contained Python programs

- Read from *standard input* and write to *standard output*!

**Mapper**

```python
#!/usr/bin/env python3

import sys
from ucb import main
from mr import emit

def emit_vowels(line):
    for vowel in 'aeiou':
        count = line.count(vowel)
        if count > 0:
            emit(vowel, count)
```

# Python Examples of a MapReduce Application

The *mapper* and *reducer* are both self-contained Python programs

- Read from *standard input* and write to *standard output*!

**Mapper**

Tell Unix: this is Python

```python
#!/usr/bin/env python3

import sys
from ucb import main
from mr import emit

def emit_vowels(line):
    for vowel in 'aeiou':
        count = line.count(vowel)
        if count > 0:
            emit(vowel, count)
```

# Python Examples of a MapReduce Application

The *mapper* and *reducer* are both self-contained Python programs

- Read from *standard input* and write to *standard output*!

**Mapper**

Tell Unix: this is Python

```python
#!/usr/bin/env python3

import sys
from ucb import main
from mr import emit

def emit_vowels(line):
    for vowel in 'aeiou':
        count = line.count(vowel)
        if count > 0:
            emit(vowel, count)
```

The emit function outputs a key and value as a line of text to standard output

# Python Examples of a MapReduce Application

The *mapper* and *reducer* are both self-contained Python programs

- Read from *standard input* and write to *standard output*!

**Mapper**

> Tell Unix: this is Python

```python
#!/usr/bin/env python3

import sys
from ucb import main
from mr import emit

def emit_vowels(line):
    for vowel in 'aeiou':
        count = line.count(vowel)
        if count > 0:
            emit(vowel, count)

@main
def run():
    for line in sys.stdin:
        emit_vowels(line)
```

> The emit function outputs a key and value as a line of text to standard output

11

# Python Examples of a MapReduce Application

The *mapper* and *reducer* are both self-contained Python programs

• Read from *standard input* and write to *standard output*!

**Mapper**

Tell Unix: this is Python

```python
#!/usr/bin/env python3

import sys
from ucb import main
from mr import emit
```

The emit function outputs a key and value as a line of text to standard output

```python
def emit_vowels(line):
    for vowel in 'aeiou':
        count = line.count(vowel)
        if count > 0:
            emit(vowel, count)

@main
def run():
    for line in sys.stdin:
        emit_vowels(line)
```

Mapper inputs are lines of text provided to standard input

The *mapper* and *reducer* are both `self-contained Python programs`

- Read from *standard input* and write to *standard output*!

**Reducer**

# Python Examples of a MapReduce Application

The *mapper* and *reducer* are both self-contained Python programs

- Read from *standard input* and write to *standard output*!

**Reducer**

```python
#!/usr/bin/env python3

import sys
from ucb import main
from mr import emit, values_by_key
```

# Python Examples of a MapReduce Application

The *mapper* and *reducer* are both self-contained Python programs

• Read from *standard input* and write to *standard output*!

**Reducer**

```python
#!/usr/bin/env python3

import sys
from ucb import main
from mr import emit, values_by_key
```

Takes and returns iterators

# Python Examples of a MapReduce Application

The *mapper* and *reducer* are both self-contained Python programs

- Read from *standard input* and write to *standard output*!

**Reducer**

```python
#!/usr/bin/env python3

import sys
from ucb import main
from mr import emit, values_by_key
```

Takes and returns iterators

**Input:** lines of text representing key-value pairs, grouped by key

**Output:** Iterator over (key, value_iterator) pairs that give all values for each key

# Python Examples of a MapReduce Application

The *mapper* and *reducer* are both self-contained Python programs

• Read from *standard input* and write to *standard output*!

**Reducer**

```python
#!/usr/bin/env python3

import sys
from ucb import main
from mr import emit, values_by_key
```

> Takes and returns iterators

> **Input:** lines of text representing key-value pairs, grouped by key
>
> **Output:** Iterator over (key, value_iterator) pairs that give all values for each key

```python
@main
def run():
    for key, value_iterator in values_by_key(sys.stdin):
        emit(key, sum(value_iterator))
```

# What Does the MapReduce Framework Provide

Tuesday, November 29, 2011

# What Does the MapReduce Framework Provide

**Fault tolerance:** A machine or hard drive might crash

# What Does the MapReduce Framework Provide

**Fault tolerance:** A machine or hard drive might crash

- The MapReduce framework automatically re-runs failed tasks.

# What Does the MapReduce Framework Provide

**Fault tolerance:** A machine or hard drive might crash

• The MapReduce framework automatically re-runs failed tasks.

**Speed:** Some machine might be slow because it's overloaded

# What Does the MapReduce Framework Provide

**Fault tolerance:** A machine or hard drive might crash

- The MapReduce framework automatically re-runs failed tasks.

**Speed:** Some machine might be slow because it's overloaded

- The framework can run multiple copies of a task and keep the result of the one that finishes first.

# What Does the MapReduce Framework Provide

**Fault tolerance:** A machine or hard drive might crash

- The MapReduce framework automatically re-runs failed tasks.

**Speed:** Some machine might be slow because it's overloaded

- The framework can run multiple copies of a task and keep the result of the one that finishes first.

**Network locality:** Data transfer is expensive

# What Does the MapReduce Framework Provide

**Fault tolerance:** A machine or hard drive might crash

- The MapReduce framework automatically re-runs failed tasks.

**Speed:** Some machine might be slow because it's overloaded

- The framework can run multiple copies of a task and keep the result of the one that finishes first.

**Network locality:** Data transfer is expensive

- The framework tries to schedule map tasks on the machines that hold the data to be processed.

Tuesday, November 29, 2011

# What Does the MapReduce Framework Provide

**Fault tolerance:** A machine or hard drive might crash

- The MapReduce framework automatically re-runs failed tasks.

**Speed:** Some machine might be slow because it's overloaded

- The framework can run multiple copies of a task and keep the result of the one that finishes first.

**Network locality:** Data transfer is expensive

- The framework tries to schedule map tasks on the machines that hold the data to be processed.

**Monitoring:** Will my job finish before dinner?!?

# What Does the MapReduce Framework Provide

**Fault tolerance:** A machine or hard drive might crash

• The MapReduce framework automatically re-runs failed tasks.

**Speed:** Some machine might be slow because it's overloaded

• The framework can run multiple copies of a task and keep the result of the one that finishes first.

**Network locality:** Data transfer is expensive

• The framework tries to schedule map tasks on the machines that hold the data to be processed.

**Monitoring:** Will my job finish before dinner?!?

• The framework provides a web-based interface describing jobs.

Tuesday, November 29, 2011