Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Regular expressions

## ComS 252 — Iowa State University

Andrew Miner

## Overview

Remember `grep`?

- ▶ "Global, regular expression, print"
- ▶ Prints lines that contain a string that matches a <span style="color:red">pattern</span>

Introduction
●○○

Simple patterns
○○○○○

Classes
○○○○○○○○

Special
○○○○○○○

Repetition
○○○○○○

Extended
○○○○○○○○

Summary
○○

## Overview

Remember grep?

▶ "Global, regular expression, print"

▶ Prints lines that contain a string that matches a pattern

How do we specify a "pattern"?

# Overview

Remember grep?

▶ "Global, regular expression, print"

▶ Prints lines that contain a string that matches a pattern

How do we specify a "pattern"?

▶ Write a regular expression

# Overview 2

- ▶ Regular expressions specify patterns to match strings
- ▶ The idea is similar to shell globbing
- ▶ Unfortunately, special characters have different meanings
  - ▶ What does globbing pattern "`hello.*`" mean?

# Overview 2

- ▶ Regular expressions specify patterns to match strings
- ▶ The idea is similar to shell globbing
- ▶ Unfortunately, special characters have different meanings
  - ▶ What does globbing pattern "`hello.*`" mean?
  - ▶ Regular expression "`hello.*`" is not the same thing

Introduction
○●○

Simple patterns
○○○○○

Classes
○○○○○○○○○

Special
○○○○○○○

Repetition
○○○○○○

Extended
○○○○○○○○○

Summary
○○

# Overview 2

▶ Regular expressions specify patterns to match strings

▶ The idea is similar to shell globbing

▶ Unfortunately, special characters have different meanings

  ▶ What does globbing pattern "`hello.*`" mean?

  ▶ Regular expression "`hello.*`" is not the same thing

    ▶ Worse — it is almost the same

# Overview 2

- ▶ Regular expressions specify patterns to match strings
- ▶ The idea is similar to shell globbing
- ▶ Unfortunately, special characters have different meanings
  - ▶ What does globbing pattern "`hello.*`" mean?
  - ▶ Regular expression "`hello.*`" is not the same thing
    - ▶ Worse — it is almost the same
- ▶ Regular expressions ("regexes") are used in many places
  - ▶ E.g., `vim`, `grep` and friends, `sed`, `awk`, `perl`
- ▶ The syntax has been standardized by POSIX
- ▶ There are two types of regular expressions:
  1. Basic regular expressions (used by `grep`)
  2. Extended regular expressions (used by `egrep`)

## Some theory

- ▶ The name comes from "regular languages"
  - ▶ Yes, the same ones from automata theory
  - ▶ See, ComS 331 is useful after all!

# Some theory

- ▶ The name comes from "regular languages"
  - ▶ Yes, the same ones from automata theory
  - ▶ See, ComS 331 is useful after all!
- ▶ We will not cover automata theory here
- ▶ Take ComS 230 / 331 for this, and learn
  - ▶ Algorithms to decide if a string matches a regex pattern
  - ▶ Why regexes include some features but not others

Introduction
○○●

Simple patterns
○○○○○

Classes
○○○○○○○○

Special
○○○○○○○

Repetition
○○○○○○

Extended
○○○○○○○○

Summary
○○

# Some theory

- The name comes from "regular languages"
    - Yes, the same ones from automata theory
    - See, ComS 331 is useful after all!
- We will not cover automata theory here
- Take ComS 230 / 331 for this, and learn
    - Algorithms to decide if a string matches a regex pattern
    - Why regexes include some features but not others
- We will discuss POSIX regex syntax

Introduction
000

Simple patterns
●0000

Classes
00000000

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

# Characters in regular expressions

There are two types of characters in a regular expression

meta characters : have special meanings

- ▶ Things like:  .   *   [   ]   \
- ▶ We will cover these, gradually

ordinary characters : everything else

Introduction
000

Simple patterns
○●○○○○

Classes
○○○○○○○○○

Special
○○○○○○○

Repetition
○○○○○○

Extended
○○○○○○○○○

Summary
○○

# Some preliminary rules

### Ordinary characters match themselves

▶ E.g., string "A" matches regex "A"

▶ E.g., string "B" does not match regex "C"

# Some preliminary rules

## Ordinary characters match themselves

- ▶ E.g., string "A" matches regex "A"
- ▶ E.g., string "B" does not match regex "C"

## Regular expressions match as much as they can, but no more

- ▶ E.g., string "Ab" does not match regex "A"
- ▶ The "as much as they can" part will make more sense later

Introduction
000

Simple patterns
000●0

Classes
00000000

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

# Concatenation rule

Concatenating regexes means concatenate strings that match

▶ More formally:
    *If* stringA *matches* regexA,
    *and* stringB *matches* regexB,
    *then* stringAstringB *matches* regexAregexB

▶ This is what you would expect

▶ E.g., string "foobar" matches regex "foobar"

▶ E.g., string "foo" does not match regex "foobar"

▶ E.g., string "bar" does not match regex "foobar"

Introduction
○○○

Simple patterns
○○○●○

Classes
○○○○○○○○○

Special
○○○○○○○

Repetition
○○○○○○

Extended
○○○○○○○○○

Summary
○○

# Concatenation rule caveat

- ▶ I stated the concatenation rule to make it look easy
- ▶ But there is a subtle catch to this rule
- ▶ Suppose I concatenate two regular expressions
    - ▶ `regexA`
    - ▶ `regexB`
- ▶ Now, when does "`string`" match `regexAregexB`?

## Concatenation rule caveat

▶ I stated the concatenation rule to make it look easy

▶ But there is a subtle catch to this rule

▶ Suppose I concatenate two regular expressions
  ▶ `regexA`
  ▶ `regexB`

▶ Now, when does "`string`" match `regexAregexB`?

▶ When there exists a way to split "`string`" such that
  1. The first "half" matches `regexA`
  2. The second "half" matches `regexB`

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

## Concatenation rule caveat

- ▶ I stated the concatenation rule to make it look easy
- ▶ But there is a subtle catch to this rule
- ▶ Suppose I concatenate two regular expressions
  - ▶ `regexA`
  - ▶ `regexB`
- ▶ Now, when does "`string`" match `regexAregexB`?
- ▶ When there exists a way to split "`string`" such that
  1. The first "half" matches `regexA`
  2. The second "half" matches `regexB`
- ▶ There are lots of ways to split a string
  - ▶ Especially when we concatenate several regexes
- ▶ You might end up matching things you do not intend
- ▶ Take care with complex regular expressions

Introduction
ooo

Simple patterns
ooooo●

Classes
oooooooo

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Example
## Carefully checking a match

String "foobar" matches regex "foobar" because
there exists a way to split string foobar:

▶ String "f" matches regex "f"

▶ String "o" matches regex "o"

▶ String "o" matches regex "o"

▶ String "b" matches regex "b"

▶ String "a" matches regex "a"

▶ String "r" matches regex "r"

Introduction
○○○

Simple patterns
○○○○○●

Classes
○○○○○○○○

Special
○○○○○○○

Repetition
○○○○○○

Extended
○○○○○○○○

Summary
○○

# Example
Carefully checking a match

String "foobar" matches regex "foobar" because
there exists a way to split string foobar:

▶ String "f" matches regex "f"

▶ String "o" matches regex "o"

▶ String "o" matches regex "o"

▶ String "b" matches regex "b"

▶ String "a" matches regex "a"

▶ String "r" matches regex "r"

This is overkill for a simple example. But it is useful when we get
to more complex rules.

# Character classes

To match a single character from a list, use [list]

► Like shell globbing
► But regexes let you do more...
► E.g., string "a" matches regex "[aeiou]"
► E.g., string "f" does not match regex "[EFL]"

Introduction
000

Simple patterns
00000

**Classes**
●00000000

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

# Character classes

To match a single character from a list, use [list]

- ▶ Like shell globbing
- ▶ But regexes let you do more...
- ▶ E.g., string "a" matches regex "[aeiou]"
- ▶ E.g., string "f" does not match regex "[EFL]"

Example

- ▶ Want to replace "folder" with "directory"
- ▶ Need to search for "folder" and "Folder"
- ▶ Can use regex:

Introduction
ooo

Simple patterns
ooooo

**Classes**
●oooooooo

Special
ooooooo

Repetition
oooooo

Extended
ooooooooo

Summary
oo

# Character classes

### To match a single character from a list, use [list]

- ▶ Like shell globbing
- ▶ But regexes let you do more...
- ▶ E.g., string "a" matches regex "[aeiou]"
- ▶ E.g., string "f" does not match regex "[EFL]"

Example

- ▶ Want to replace "folder" with "directory"
- ▶ Need to search for "folder" and "Folder"
- ▶ Can use regex:
  [Ff]older

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

# Ranges in character classes

## Use "–" inside [] to specify a character range

- ▶ E.g., use regex "[01]" to match a binary digit
- ▶ E.g., use regex "[0-7]" to match an octal digit
- ▶ E.g., use regex "[0-9]" to match a decimal digit
- ▶ E.g., use regex "[0-9a-f]" to match a hexadecimal digit

Introduction
ooo

Simple patterns
ooooo

**Classes**
oooooooo

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Ranges in character classes

## Use "–" inside [] to specify a character range

- ▶ E.g., use regex "[01]" to match a binary digit
- ▶ E.g., use regex "[0-7]" to match an octal digit
- ▶ E.g., use regex "[0-9]" to match a decimal digit
- ▶ E.g., use regex "[0-9a-f]" to match a hexadecimal digit

## Use "–" first or last inside [] to specify "–"

- ▶ E.g., to match arithmetic operators, use "[-+*/]"

Introduction
000

Simple patterns
00000

**Classes**
00●00000

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

# Example

How can we match *all* text of the form:

1. Single–character variable name
2. Space
3. Arithmetic operator
4. Space
5. Single–character variable name

Things like: $x + y$, $I - J$, ...

# Example

How can we match *all* text of the form:

1. Single–character variable name
2. Space
3. Arithmetic operator
4. Space
5. Single–character variable name

Things like: $x + y$, $I - J$, . . .
Use regex: "[a-zA-Z] [-+*/%] [a-zA-Z]"

Introduction
ooo

Simple patterns
ooooo

**Classes**
ooo●oooo

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Testing these things

The easiest way is to send a string through grep:

```
prompt$
```

Introduction
000

Simple patterns
00000

**Classes**
000●0000

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

# Testing these things

The easiest way is to send a string through grep:

```
prompt$ echo "folder" | grep "[fF]older"
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooooo

Special
ooooooo

Repetition
oooooo

Extended
oooooooooo

Summary
oo

# Testing these things

The easiest way is to send a string through grep:

```
prompt$ echo "folder" | grep "[fF]older"
folder
prompt$ 
```

Introduction
ooo

Simple patterns
ooooo

**Classes**
ooo●oooo

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Testing these things

The easiest way is to send a string through grep:

```
prompt$ echo "folder" | grep "[fF]older"
folder
prompt$ echo "[fF]older" | grep "[fF]older"
```

Introduction
000

Simple patterns
00000

**Classes**
00000000

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

# Testing these things

The easiest way is to send a string through grep:

```
prompt$ echo "folder" | grep "[fF]older"
folder
prompt$ echo "[fF]older" | grep "[fF]older"
prompt$
```

Introduction
ooo

Simple patterns
ooooo

**Classes**
ooooooooo

Special
ooooooo

Repetition
oooooo

Extended
ooooooooo

Summary
oo

# Testing these things

The easiest way is to send a string through grep:

```
prompt$ echo "folder" | grep "[fF]older"
folder
prompt$ echo "[fF]older" | grep "[fF]older"
prompt$ echo "x + y" | grep "[a-z] [-+] [a-z]"
```

Introduction
000

Simple patterns
00000

**Classes**
00000000

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

# Testing these things

The easiest way is to send a string through grep:

```
prompt$ echo "folder" | grep "[fF]older"
folder
prompt$ echo "[fF]older" | grep "[fF]older"
prompt$ echo "x + y" | grep "[a-z] [-+] [a-z]"
x + y
prompt$ ▮
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooooo

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Testing these things

The easiest way is to send a string through grep:

```
prompt$ echo "folder" | grep "[fF]older"
folder
prompt$ echo "[fF]older" | grep "[fF]older"
prompt$ echo "x + y" | grep "[a-z] [-+] [a-z]"
x + y
prompt$ echo $?
```

Introduction
ooo

Simple patterns
ooooo

**Classes**
oooo●oooo

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Testing these things

The easiest way is to send a string through grep:

```
prompt$ echo "folder" | grep "[fF]older"
folder
prompt$ echo "[fF]older" | grep "[fF]older"
prompt$ echo "x + y" | grep "[a-z] [-+] [a-z]"
x + y
prompt$ echo $?
0
prompt$ █
```

Introduction
ooo

Simple patterns
ooooo

**Classes**
ooo●oooo

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

## Testing these things

The easiest way is to send a string through grep:

```
prompt$ echo "folder" | grep "[fF]older"
folder
prompt$ echo "[fF]older" | grep "[fF]older"
prompt$ echo "x + y" | grep "[a-z] [-+] [a-z]"
x + y
prompt$ echo $?
0
prompt$ echo "x+y" | grep "[a-z] [-+] [a-z]"
```

# Testing these things

The easiest way is to send a string through grep:

```
prompt$ echo "folder" | grep "[fF]older"
folder
prompt$ echo "[fF]older" | grep "[fF]older"
prompt$ echo "x + y" | grep "[a-z] [-+] [a-z]"
x + y
prompt$ echo $?
0
prompt$ echo "x+y" | grep "[a-z] [-+] [a-z]"
prompt$
```

Introduction
000

Simple patterns
00000

**Classes**
00000000

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

# Testing these things

The easiest way is to send a string through grep:

```
prompt$ echo "folder" | grep "[fF]older"
folder
prompt$ echo "[fF]older" | grep "[fF]older"
prompt$ echo "x + y" | grep "[a-z] [-+] [a-z]"
x + y
prompt$ echo $?
0
prompt$ echo "x+y" | grep "[a-z] [-+] [a-z]"
prompt$ echo $?
```

Introduction
000

Simple patterns
00000

**Classes**
00000000

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

## Testing these things

The easiest way is to send a string through grep:

```
prompt$ echo "folder" | grep "[fF]older"
folder
prompt$ echo "[fF]older" | grep "[fF]older"
prompt$ echo "x + y" | grep "[a-z] [-+] [a-z]"
x + y
prompt$ echo $?
0
prompt$ echo "x+y" | grep "[a-z] [-+] [a-z]"
prompt$ echo $?
1
prompt$
```

grep has a zero exit code if there was at least one matching line

Introduction
ooo

Simple patterns
ooooo

**Classes**
ooooo●ooo

Special
ooooooo

Repetition
oooooo

Extended
ooooooooo

Summary
oo

# Inverting a class

### Use "^" as the first character in [] to invert the class

▶ E.g., use regex "[^f]" to match every character except f

▶ E.g., use regex "[^aeiou]" to match consonants

# Inverting a class

## Use "^" as the first character in [] to invert the class

▶ E.g., use regex "[^f]" to match every character except f

▶ E.g., use regex "[^aeiou]" to match consonants
and digits, and control characters, and many other things
but not lower–case vowels

Introduction
ooo

Simple patterns
ooooo

**Classes**
oooooo●oo

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Tricky ones

What do the following regular expressions match?

`[^^]`

# Tricky ones

What do the following regular expressions match?

[^^]

▶ Any character except "^"

Introduction
000

Simple patterns
00000

**Classes**
00000●00

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

# Tricky ones

What do the following regular expressions match?

[^^]

▶ Any character except "^"

[^-z]

Introduction
ooo

Simple patterns
ooooo

Classes
oooooo●oo

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Tricky ones

What do the following regular expressions match?

[^^]

▶ Any character except "^"

[^-z]

▶ Any character between "^" and "z"
▶ Any character except "-" and "z"

Introduction
ooo

Simple patterns
ooooo

**Classes**
ooooooo●oo

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Tricky ones

What do the following regular expressions match?

     [^^]

> ▶ Any character except "^"

     [^-z]

> ▶ ~~Any character between "^" and "z"~~ NO
> ▶ Any character except "-" and "z"

Introduction
ooo

Simple patterns
ooooo

Classes
ooooooo●oo

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

## Tricky ones

What do the following regular expressions match?

      `[^^]`

> ▶ Any character except "^"

      `[^-z]`

> ▶ ~~Any character between "^" and "z"~~ NO
> ▶ Any character except "-" and "z"

How do we write: any character between "^" and "z"?

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○●○○

Special
○○○○○○○

Repetition
○○○○○○

Extended
○○○○○○○○

Summary
○○

# Tricky ones

What do the following regular expressions match?

[^^]

    ▶ Any character except "^"

[^-z]

    ▶ ~~Any character between "^" and "z"~~ NO
    ▶ Any character except "-" and "z"

How do we write: any character between "^" and "z"?

"[z^-z]"

Introduction
000

Simple patterns
00000

**Classes**
000000●00

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

# Tricky ones

What do the following regular expressions match?

    `[^^]`

        ► Any character except "^"

    `[^-z]`

        ► ~~Any character between "^" and "z"~~ NO
        ► Any character except "-" and "z"

How do we write: any character between "^" and "z"?

    "`[z^-z]`"

What is a *safer* way to write: any character except "-" and "z"?

Introduction
○○○

Simple patterns
○○○○○

**Classes**
○○○○○○●○○

Special
○○○○○○○

Repetition
○○○○○○

Extended
○○○○○○○○

Summary
○○

# Tricky ones

What do the following regular expressions match?

      `[^^]`

> ▶ Any character except "^"

      `[^-z]`

> ▶ ~~Any character between "^" and "z"~~ NO
> ▶ Any character except "−" and "z"

How do we write: any character between "^" and "z"?

      "`[z^-z]`"

What is a *safer* way to write: any character except "−" and "z"?

      "`[^z-]`"

Introduction
ooo

Simple patterns
ooooo

**Classes**
ooooooo●o

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Names for classes

POSIX defines names for classes of characters, including:

[:alpha:] : Alphabetic characters

Introduction
ooo

Simple patterns
ooooo

**Classes**
ooooooo●o

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

## Names for classes

POSIX defines names for classes of characters, including:

[:alpha:] : Alphabetic characters

▶ Isn't this the same as "a-zA-Z"?

Introduction
ooo

Simple patterns
ooooo

**Classes**
ooooooo●o

Special
ooooooo

Repetition
oooooo

Extended
ooooooooo

Summary
oo

## Names for classes

POSIX defines names for classes of characters, including:

[:alpha:] : Alphabetic characters

> ▶ Isn't this the same as "a-zA-Z"?
> ▶ Not always: [:alpha:] may match "è"

# Names for classes

POSIX defines names for classes of characters, including:

[:alpha:] : Alphabetic characters

  ▶ Isn't this the same as "a-zA-Z"?
  ▶ Not always: [:alpha:] may match "è"

[:alnum:] : Alphanumeric characters

[:cntrl:] : Control characters

[:digit:] : Numeric characters

[:lower:] : Lower–case characters

[:punct:] : Punctuation characters

[:space:] : Space, tab, and other "whitespace"

[:upper:] : Upper–case characters

Introduction
○○○

Simple patterns
○○○○○

**Classes**
○○○○○○○●○

Special
○○○○○○○

Repetition
○○○○○○

Extended
○○○○○○○○

Summary
○○

## Names for classes

POSIX defines names for classes of characters, including:

[:alpha:] : Alphabetic characters
- ▶ Isn't this the same as "a-zA-Z"?
- ▶ Not always: [:alpha:] may match "è"

[:alnum:] : Alphanumeric characters

[:cntrl:] : Control characters

[:digit:] : Numeric characters

[:lower:] : Lower–case characters

[:punct:] : Punctuation characters

[:space:] : Space, tab, and other "whitespace"

[:upper:] : Upper–case characters

These must be used within []

Introduction
ooo

Simple patterns
ooooo

Classes
ooooooo●

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Example

```
prompt$ █
```

# Example

```
prompt$ echo "a" | grep "[:alpha:]"
```

Introduction
ooo

Simple patterns
ooooo

Classes
ooooooo●

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Example

```
prompt$ echo "a" | grep "[:alpha:]"
a
prompt$
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo●

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Example

```
prompt$ echo "a" | grep "[:alpha:]"
a
prompt$ echo "b" | grep "[:alpha:]"
```

Introduction
ooo

Simple patterns
ooooo

**Classes**
ooooooo●

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Example

```
prompt$ echo "a" | grep "[:alpha:]"
a
prompt$ echo "b" | grep "[:alpha:]"
prompt$ █
```

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○●

Special
○○○○○○○

Repetition
○○○○○○

Extended
○○○○○○○○

Summary
○○

# Example

```
prompt$ echo "a" | grep "[:alpha:]"
a
prompt$ echo "b" | grep "[:alpha:]"
prompt$ echo "a" | grep "[[:alpha:]]"
```

Introduction
ooo
Simple patterns
ooooo
Classes
ooooooo●
Special
ooooooo
Repetition
oooooo
Extended
ooooooooo
Summary
oo

# Example

```
prompt$ echo "a" | grep "[:alpha:]"
a
prompt$ echo "b" | grep "[:alpha:]"
prompt$ echo "a" | grep "[[:alpha:]]"
a
prompt$ ▮
```

Introduction
ooo

Simple patterns
ooooo

**Classes**
ooooooo●

Special
ooooooo

Repetition
oooooo

Extended
ooooooooo

Summary
oo

# Example

```
prompt$ echo "a" | grep "[:alpha:]"
a
prompt$ echo "b" | grep "[:alpha:]"
prompt$ echo "a" | grep "[[:alpha:]]"
a
prompt$ echo "b" | grep "[[:alpha:]]"
```

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○●

Special
○○○○○○○

Repetition
○○○○○○

Extended
○○○○○○○○

Summary
○○

# Example

```
prompt$ echo "a" | grep "[:alpha:]"
a
prompt$ echo "b" | grep "[:alpha:]"
prompt$ echo "a" | grep "[[:alpha:]]"
a
prompt$ echo "b" | grep "[[:alpha:]]"
b
prompt$ 
```

Introduction
ooo
Simple patterns
ooooo
Classes
ooooooo●
Special
ooooooo
Repetition
oooooo
Extended
oooooooo
Summary
oo

# Example

```
prompt$ echo "a" | grep "[:alpha:]"
a
prompt$ echo "b" | grep "[:alpha:]"
prompt$ echo "a" | grep "[[:alpha:]]"
a
prompt$ echo "b" | grep "[[:alpha:]]"
b
prompt$ echo "ò" | grep "[[:alpha:]]"
```

Introduction
000

Simple patterns
00000

**Classes**
0000000●

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

# Example

```
prompt$ echo "a" | grep "[:alpha:]"
a
prompt$ echo "b" | grep "[:alpha:]"
prompt$ echo "a" | grep "[[:alpha:]]"
a
prompt$ echo "b" | grep "[[:alpha:]]"
b
prompt$ echo "ò" | grep "[[:alpha:]]"
ò
prompt$ █
```

I tested this on a Macintosh, and used `Option-`'o` to get "ò"

Introduction
000

Simple patterns
00000

**Classes**
0000000●

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

# Example

```
prompt$ echo "a" | grep "[:alpha:]"
a
prompt$ echo "b" | grep "[:alpha:]"
prompt$ echo "a" | grep "[[:alpha:]]"
a
prompt$ echo "b" | grep "[[:alpha:]]"
b
prompt$ echo "ò" | grep "[[:alpha:]]"
ò
prompt$ echo "9" | grep "[[:alpha:]]"
```

I tested this on a Macintosh, and used Option-'o to get "ò"

Introduction
ooo

Simple patterns
ooooo

**Classes**
oooooooo●

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Example

```
prompt$ echo "a" | grep "[:alpha:]"
a
prompt$ echo "b" | grep "[:alpha:]"
prompt$ echo "a" | grep "[[:alpha:]]"
a
prompt$ echo "b" | grep "[[:alpha:]]"
b
prompt$ echo "ò" | grep "[[:alpha:]]"
ò
prompt$ echo "9" | grep "[[:alpha:]]"
prompt$ █
```

I tested this on a Macintosh, and used `Option-'o` to get "ò"

Introduction
000

Simple patterns
00000

Classes
0000000●

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

# Example

```
prompt$ echo "a" | grep "[:alpha:]"
a
prompt$ echo "b" | grep "[:alpha:]"
prompt$ echo "a" | grep "[[:alpha:]]"
a
prompt$ echo "b" | grep "[[:alpha:]]"
b
prompt$ echo "ò" | grep "[[:alpha:]]"
ò
prompt$ echo "9" | grep "[[:alpha:]]"
prompt$ echo "9" | grep "[^[:alpha:]]"
```

I tested this on a Macintosh, and used Option-`o to get "ò"

Introduction
ooo

Simple patterns
ooooo

**Classes**
ooooooo●

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Example

```
prompt$ echo "a" | grep "[:alpha:]"
a
prompt$ echo "b" | grep "[:alpha:]"
prompt$ echo "a" | grep "[[:alpha:]]"
a
prompt$ echo "b" | grep "[[:alpha:]]"
b
prompt$ echo "ò" | grep "[[:alpha:]]"
ò
prompt$ echo "9" | grep "[[:alpha:]]"
prompt$ echo "9" | grep "[^[:alpha:]]"
9
prompt$ █
```

I tested this on a Macintosh, and used `Option-'o` to get "ò"

# Some special characters

### The "." character matches (almost) any single character

▶ The newline character usually does not match "."

▶ "." would be equivalent to "[^<newline>]"
   if you could specify the newline character in a class

```
prompt$ 
```

Introduction
000
Simple patterns
00000
Classes
00000000
Special
●000000
Repetition
000000
Extended
00000000
Summary
00

# Some special characters

## The "." character matches (almost) any single character

▶ The newline character usually does not match "."

▶ "." would be equivalent to "[^<newline>]"

  if you could specify the newline character in a class

```
prompt$ echo "ab" | grep "a."
```

## Some special characters

The "." character matches (almost) any single character

▶ The newline character usually does not match "."

▶ "." would be equivalent to "[^<newline>]"

   if you could specify the newline character in a class

```
prompt$ echo "ab" | grep "a."
ab
prompt$ 
```

# Some special characters

### The "." character matches (almost) any single character

► The newline character usually does not match "."

► "." would be equivalent to "[^<newline>]"

  if you could specify the newline character in a class

```
prompt$ echo "ab" | grep "a."
ab
prompt$ echo "a" | grep "a."
```

# Some special characters

### The "." character matches (almost) any single character

▶ The newline character usually does not match "."

▶ "." would be equivalent to "[^<newline>]"

   if you could specify the newline character in a class

```
prompt$ echo "ab" | grep "a."
ab
prompt$ echo "a" | grep "a."
prompt$
```

Introduction
ooo
Simple patterns
ooooo
Classes
oooooooo
Special
o●ooooo
Repetition
oooooo
Extended
oooooooo
Summary
oo

How to match an actual "." (only)?

Introduction
ooo
Simple patterns
ooooo
Classes
oooooooo
Special
o●ooooo
Repetition
oooooo
Extended
oooooooo
Summary
oo

# How to match an actual "." (only)?

1. Use "[.]"

# How to match an actual "." (only)?

1. Use "[.]"

```
prompt$ ▋
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

**Special**
oooooooo

Repetition
oooooo

Extended
ooooooooo

Summary
oo

# How to match an actual "." (only)?

1. Use "[.]"

```
prompt$ echo "ab" | grep "a[.]"
```

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○○

**Special**
○●○○○○○

Repetition
○○○○○○

Extended
○○○○○○○○○

Summary
○○

# How to match an actual "." (only)?

1. Use "[.]"

```
prompt$ echo "ab" | grep "a[.]"
prompt$
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

**Special**
o●ooooo

Repetition
oooooo

Extended
ooooooooo

Summary
oo

# How to match an actual "." (only)?

1. Use "[.]"

```
prompt$ echo "ab" | grep "a[.]"
prompt$ echo "a." | grep "a[.]"
```

# How to match an actual "." (only)?

1. Use "[.]"

```
prompt$ echo "ab" | grep "a[.]"
prompt$ echo "a." | grep "a[.]"
a.
prompt$ █
```

Introduction
ooo
Simple patterns
ooooo
Classes
oooooooo
**Special**
oooooooo
Repetition
oooooo
Extended
ooooooooo
Summary
oo

# How to match an actual "." (only)?

1. Use "[.]"

```
prompt$ echo "ab" | grep "a[.]"
prompt$ echo "a." | grep "a[.]"
a.
prompt$
```

2. Use "\."

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

**Special**
oooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# How to match an actual "." (only)?

1. Use "[.]"

```
prompt$ echo "ab" | grep "a[.]"
prompt$ echo "a." | grep "a[.]"
a.
prompt$
```

2. Use "\."

```
prompt$
```

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○○

**Special**
○●○○○○○

Repetition
○○○○○○

Extended
○○○○○○○○

Summary
○○

# How to match an actual "." (only)?

1. Use "[.]"

```
prompt$ echo "ab" | grep "a[.]"
prompt$ echo "a." | grep "a[.]"
a.
prompt$
```

2. Use "\."

```
prompt$ echo "ab" | grep "a\."
```

# How to match an actual "." (only)?

1. Use "[.]"

```
prompt$ echo "ab" | grep "a[.]"
prompt$ echo "a." | grep "a[.]"
a.
prompt$
```

2. Use "\."

```
prompt$ echo "ab" | grep "a\."
prompt$
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
o●ooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# How to match an actual "." (only)?

1. Use "[.]"

```
prompt$ echo "ab" | grep "a[.]"
prompt$ echo "a." | grep "a[.]"
a.
prompt$
```

2. Use "\."

```
prompt$ echo "ab" | grep "a\."
prompt$ echo "a." | grep "a\."
```

# How to match an actual "." (only)?

1. Use "[.]"

```
prompt$ echo "ab" | grep "a[.]"
prompt$ echo "a." | grep "a[.]"
a.
prompt$
```

2. Use "\."

```
prompt$ echo "ab" | grep "a\."
prompt$ echo "a." | grep "a\."
a.
prompt$ █
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Escaping characters

### The "\" character will "escape" the following character

▶ Converts meta characters into ordinary ones

▶ Be careful because this is also a special shell character

Example: match the string "["

```
prompt$ ▊
```

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○○

**Special**
○○●○○○○○

Repetition
○○○○○○

Extended
○○○○○○○○

Summary
○○

## Escaping characters

The "\" character will "escape" the following character

▶ Converts meta characters into ordinary ones
▶ Be careful because this is also a special shell character

Example: match the string "["

```
prompt$ echo '['
```

Introduction
ooo
Simple patterns
ooooo
Classes
oooooooo
**Special**
ooo●oooo
Repetition
oooooo
Extended
oooooooo
Summary
oo

# Escaping characters

The "\" character will "escape" the following character

- ▶ Converts meta characters into ordinary ones
- ▶ Be careful because this is also a special shell character

Example: match the string "["

```
prompt$ echo '['
[
prompt$ █
```

Introduction
000

Simple patterns
00000

Classes
00000000

**Special**
0000000

Repetition
000000

Extended
00000000

Summary
00

# Escaping characters

The "\" character will "escape" the following character

▶ Converts meta characters into ordinary ones
▶ Be careful because this is also a special shell character

Example: match the string "["

```
prompt$ echo '['
[
prompt$ echo '[' | grep '['
```

# Escaping characters

The "\" character will "escape" the following character

▶ Converts meta characters into ordinary ones
▶ Be careful because this is also a special shell character

Example: match the string "["

```
prompt$ echo '['
[
prompt$ echo '[' | grep '['
grep: Unmatched [ or [^
prompt$
```

# Escaping characters

## The "\" character will "escape" the following character

▶ Converts meta characters into ordinary ones
▶ Be careful because this is also a special shell character

Example: match the string "["

```
prompt$ echo '['
[
prompt$ echo '[' | grep '['
grep: Unmatched [ or [^
prompt$ echo '[' | grep '\['
```

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○○

**Special**
○○○●○○○○

Repetition
○○○○○○

Extended
○○○○○○○○

Summary
○○

# Escaping characters

The "\" character will "escape" the following character

▶ Converts meta characters into ordinary ones
▶ Be careful because this is also a special shell character

Example: match the string "["

```
prompt$ echo '['
[
prompt$ echo '[' | grep '['
grep: Unmatched [ or [^
prompt$ echo '[' | grep '\['
[
prompt$ 
```

Introduction
000

Simple patterns
00000

Classes
00000000

**Special**
0000●000

Repetition
000000

Extended
00000000

Summary
00

# Example: match the string "\"

```
prompt$
```

Introduction
ooo
Simple patterns
ooooo
Classes
oooooooo
**Special**
oooo●ooo
Repetition
oooooo
Extended
oooooooo
Summary
oo

# Example: match the string "\"

```
prompt$ echo \
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

**Special**
oooo●ooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Example: match the string "\"

```
prompt$ echo \
>
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

**Special**
ooooo●ooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Example: match the string "\"

```
prompt$ echo \
> What is happening here
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

**Special**
ooo●ooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Example: match the string "\"

```
prompt$ echo \
> What is happening here
What is happening here
prompt$
```

# Example: match the string "\"

```
prompt$ echo \
> What is happening here
What is happening here
prompt$ █
```

▶ \ at the end of a line means "let me continue on the next line"

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

**Special**
ooo●ooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Example: match the string "\"

```
prompt$ echo \
> What is happening here
What is happening here
prompt$ echo "\"
```

▶ \ at the end of a line means "let me continue on the next line"

# Example: match the string "\"

```
prompt$ echo \
> What is happening here
What is happening here
prompt$ echo "\"
>
```

▶ \ at the end of a line means "let me continue on the next line"

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

**Special**
ooooooo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Example: match the string "\"

```
prompt$ echo \
> What is happening here
What is happening here
prompt$ echo "\"
> Again?
```

▶ \ at the end of a line means "let me continue on the next line"

Introduction
000

Simple patterns
00000

Classes
00000000

**Special**
0000●000

Repetition
000000

Extended
00000000

Summary
00

# Example: match the string "\"

```
prompt$ echo \
> What is happening here
What is happening here
prompt$ echo "\"
> Again?
>
```

▶ \ at the end of a line means "let me continue on the next line"

Introduction
ooo
Simple patterns
ooooo
Classes
oooooooo
**Special**
ooooo●ooo
Repetition
oooooo
Extended
oooooooo
Summary
oo

# Example: match the string "\"

```
prompt$ echo \
> What is happening here
What is happening here
prompt$ echo "\"
> Again?
> "
```

▶ \ at the end of a line means "let me continue on the next line"

Introduction
000

Simple patterns
00000

Classes
00000000

**Special**
0000●000

Repetition
000000

Extended
00000000

Summary
00

# Example: match the string "\"

```
prompt$ echo \
> What is happening here
What is happening here
prompt$ echo "\"
> Again?
> "
"
Again?

prompt$ █
```

▶ \ at the end of a line means "let me continue on the next line"

Introduction
000

Simple patterns
00000

Classes
00000000

**Special**
0000●000

Repetition
000000

Extended
00000000

Summary
00

# Example: match the string "\"

```
prompt$ echo \
> What is happening here
What is happening here
prompt$ echo "\"
> Again?
> "
"
Again?

prompt$ █
```

▶ \ at the end of a line means "let me continue on the next line"
▶ \" converts " to a normal shell character

# Example: match the string "\"

```
prompt$ echo \
> What is happening here
What is happening here
prompt$ echo "\"
> Again?
> "
"
Again?

prompt$ echo \"
```

► \ at the end of a line means "let me continue on the next line"
► \" converts " to a normal shell character

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000●000

Repetition
000000

Extended
00000000

Summary
00

# Example: match the string "\"

```
What is happening here
prompt$ echo "\"
> Again?
> "
"
Again?

prompt$ echo \"
"
prompt$ █
```

▶ \ at the end of a line means "let me continue on the next line"
▶ \" converts " to a normal shell character

# Example: match the string "\"

```
What is happening here
prompt$ echo "\"
> Again?
> "
"
Again?

prompt$ echo \"
"
prompt$ echo "\\"█
```

▶ \ at the end of a line means "let me continue on the next line"
▶ \" converts " to a normal shell character

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

**Special**
ooooo●oo

Repetition
oooooo

Extended
ooooooooo

Summary
oo

# Example: match the string "\"

```
> Again?
> "
"
Again?

prompt$ echo \"
"
prompt$ echo "\\"
\
prompt$ ▮
```

- ▶ \ at the end of a line means "let me continue on the next line"
- ▶ \" converts " to a normal shell character
- ▶ \\ converts \ to a normal shell character

Introduction
000

Simple patterns
00000

Classes
00000000

**Special**
0000●000

Repetition
000000

Extended
00000000

Summary
00

# Example: match the string "\"

```
> Again?
> "
"
Again?

prompt$ echo \"
"
prompt$ echo "\\"
\
prompt$ echo '\'
```

- ▶ \ at the end of a line means "let me continue on the next line"
- ▶ \" converts " to a normal shell character
- ▶ \\ converts \ to a normal shell character

Introduction
000

Simple patterns
00000

Classes
00000000

**Special**
0000●000

Repetition
000000

Extended
00000000

Summary
00

# Example: match the string "\"

```
"
Again?

prompt$ echo \"
"
prompt$ echo "\\"
\
prompt$ echo '\'
\
prompt$ 
```

- ► \ at the end of a line means "let me continue on the next line"
- ► \" converts " to a normal shell character
- ► \\ converts \ to a normal shell character

Introduction
000

Simple patterns
00000

Classes
00000000

**Special**
0000●000

Repetition
000000

Extended
00000000

Summary
00

# Example: match the string "\"

```
"
Again?

prompt$ echo \"
"
prompt$ echo "\\"
\
prompt$ echo '\'
\
prompt$ echo '\' | grep '\'
```

- ▶ \ at the end of a line means "let me continue on the next line"
- ▶ \" converts " to a normal shell character
- ▶ \\ converts \ to a normal shell character

Introduction
000

Simple patterns
00000

Classes
00000000

**Special**
0000●000

Repetition
000000

Extended
00000000

Summary
00

# Example: match the string "\"

```
prompt$ echo \"
"
prompt$ echo "\\"
\
prompt$ echo '\'
\
prompt$ echo '\' | grep '\'
grep: Trailing backslash
prompt$ ▊
```

- ▶ \ at the end of a line means "let me continue on the next line"
- ▶ \" converts " to a normal shell character
- ▶ \\ converts \ to a normal shell character

Introduction
000

Simple patterns
00000

Classes
00000000

**Special**
0000●000

Repetition
000000

Extended
00000000

Summary
00

# Example: match the string "\"

```
prompt$ echo \"
"
prompt$ echo "\\"
\
prompt$ echo '\'
\
prompt$ echo '\' | grep '\'
grep: Trailing backslash
prompt$ echo '\' | grep '\\'
```

- ▶ \ at the end of a line means "let me continue on the next line"
- ▶ \" converts " to a normal shell character
- ▶ \\ converts \ to a normal shell character

Introduction
000

Simple patterns
00000

Classes
00000000

**Special**
0000●000

Repetition
000000

Extended
00000000

Summary
00

# Example: match the string "\"

```
"
prompt$ echo "\\"
\
prompt$ echo '\'
\
prompt$ echo '\' | grep '\'
grep: Trailing backslash
prompt$ echo '\' | grep '\\'
\
prompt$ ▋
```

► \ at the end of a line means "let me continue on the next line"

► \" converts " to a normal shell character

► \\ converts \ to a normal shell character

Introduction
000

Simple patterns
00000

Classes
00000000

**Special**
0000000

Repetition
000000

Extended
00000000

Summary
00

# Example: match the string "\"

```
"
prompt$ echo "\\"
\
prompt$ echo '\'
\
prompt$ echo '\' | grep '\'
grep: Trailing backslash
prompt$ echo '\' | grep '\\'
\
prompt$ echo '\' | grep "\\\\"
```

- ▶ \ at the end of a line means "let me continue on the next line"
- ▶ \" converts " to a normal shell character
- ▶ \\ converts \ to a normal shell character

Introduction
000

Simple patterns
00000

Classes
00000000

**Special**
0000●000

Repetition
000000

Extended
00000000

Summary
00

# Example: match the string "\"

```
\
prompt$ echo '\'
\
prompt$ echo '\' | grep '\'
grep: Trailing backslash
prompt$ echo '\' | grep '\\'
\
prompt$ echo '\' | grep "\\\\"
\
prompt$ █
```

▶ \ at the end of a line means "let me continue on the next line"

▶ \" converts " to a normal shell character

▶ \\ converts \ to a normal shell character

▶ Probably best to use single quotes for regular expressions

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

**Special**
oooo●oo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Motivating example

```
prompt$ 
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

**Special**
oooo●oo

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Motivating example

```
prompt$ echo "test" | grep 'e'
```

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○○○

**Special**
○○○○●○○

Repetition
○○○○○○

Extended
○○○○○○○○○

Summary
○○

## Motivating example

```
prompt$ echo "test" | grep 'e'
test
prompt$ █
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

**Special**
oooo●oo

Repetition
oooooo

Extended
ooooooooo

Summary
oo

# Motivating example

```
prompt$ echo "test" | grep 'e'
test
prompt$
```

What just happened?

# Motivating example

```
prompt$ echo "test" | grep 'e'
test
prompt$
```

What just happened?

- "e" matches regular expression 'e'

# Motivating example

```
prompt$ echo "test" | grep 'e'
test
prompt$
```

What just happened?

- ▶ "e" matches regular expression 'e'
- ▶ The line "test" contains text that matches the pattern
    - ▶ So the whole line is printed
    - ▶ Remember how grep works?
    - ▶ Print lines containing text that matches a pattern

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○○

**Special**
○○○○●○○

Repetition
○○○○○○

Extended
○○○○○○○○

Summary
○○

## Motivating example

```
prompt$ echo "test" | grep 'e'
test
prompt$
```

What just happened?

- ▶ "e" matches regular expression 'e'
- ▶ The line "test" contains text that matches the pattern
  - ▶ So the whole line is printed
  - ▶ Remember how grep works?
  - ▶ Print lines containing text that matches a pattern

How can I force grep to match only the line 'e'?

Introduction
ooo
Simple patterns
ooooo
Classes
oooooooo
Special
oooooo●o
Repetition
oooooo
Extended
ooooooooo
Summary
oo

# More special characters

## Character "^" matches "beginning of line"

▶ There is no "beginning of line" character

▶ It is a special, "imaginary character"

## Character "$" matches "end of line"

▶ Not the newline character

▶ It is a special, "imaginary character"

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

**Special**
ooooooeo

Repetition
oooooo

Extended
ooooooooo

Summary
oo

## More special characters

### Character "^" matches "beginning of line"

▶ There is no "beginning of line" character

▶ It is a special, "imaginary character"

### Character "$" matches "end of line"

▶ Not the newline character

▶ It is a special, "imaginary character"

Fun fact: guess what "^" and "$" do in `vi`?

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○○

Special
○○○○○●○

Repetition
○○○○○○

Extended
○○○○○○○○

Summary
○○

# More special characters

## Character "^" matches "beginning of line"

▶ There is no "beginning of line" character

▶ It is a special, "imaginary character"

## Character "$" matches "end of line"

▶ Not the newline character

▶ It is a special, "imaginary character"

Fun fact: guess what "^" and "$" do in vi?

    ^ moves the cursor to the beginning of the line

    $ moves the cursor to the end of the line

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

**Special**
oooooo●

Repetition
oooooo

Extended
ooooooooo

Summary
oo

# Examples

▶ To match lines that <span style="color:red">start</span> with "a":

Introduction
ooo
Simple patterns
ooooo
Classes
oooooooo
**Special**
oooooo●
Repetition
oooooo
Extended
ooooooooo
Summary
oo

# Examples

- To match lines that <span style="color:red">start</span> with "a":

  grep '^a'

# Examples

▶ To match lines that start with "a":

  grep '^a'

▶ To match lines that end with "a":

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

**Special**
oooooo●

Repetition
oooooo

Extended
ooooooooo

Summary
oo

# Examples

- ▶ To match lines that <span style="color:red">start</span> with "a":

  `grep '^a'`

- ▶ To match lines that <span style="color:red">end</span> with "a":

  `grep 'a$'`

# Examples

- To match lines that start with "a":

  grep '^a'
- To match lines that end with "a":

  grep 'a$'
- To match the line "a":

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo●

Repetition
oooooo

Extended
ooooooooo

Summary
oo

# Examples

▶ To match lines that start with "a":

`grep '^a'`

▶ To match lines that end with "a":

`grep 'a$'`

▶ To match the line "a":

`grep '^a$'`

# Examples

▶ To match lines that <span style="color:red">start</span> with "a":

  `grep '^a'`

▶ To match lines that <span style="color:red">end</span> with "a":

  `grep 'a$'`

▶ To match the line "a":

  `grep '^a$'`

▶ To match a line containing "^":

# Examples

- To match lines that start with "a":

  grep '^a'
- To match lines that end with "a":

  grep 'a$'
- To match the line "a":

  grep '^a$'
- To match a line containing "^":

  grep '\^'

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
oooooo●

Repetition
oooooo

Extended
oooooooo

Summary
oo

# Examples

- To match lines that start with "a":

  grep '^a'
- To match lines that end with "a":

  grep 'a$'
- To match the line "a":

  grep '^a$'
- To match a line containing "^":

  grep '\^'
- To match a line containing "$":

# Examples

- ▶ To match lines that <span style="color:red">start</span> with "a":

  grep '^a'

- ▶ To match lines that <span style="color:red">end</span> with "a":

  grep 'a$'

- ▶ To match the line "a":

  grep '^a$'

- ▶ To match a line containing "^":

  grep '\^'

- ▶ To match a line containing "$":

  grep '\$'

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
●ooooo

Extended
oooooooo

Summary
oo

# Repeating things

This is where regular expressions get tricky to match, by hand

# Repeating things

This is where regular expressions get tricky to match, by hand

"∗" means, repeat the previous thing, zero or more times

▶ The "previous thing" is a character
  ▶ For basic regular expressions, anyway

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
o●oooo

Extended
oooooooo

Summary
oo

# Simple example with "∗"

▶ Does string 'bd' match regex 'ba∗d'?

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
o●oooo

Extended
oooooooo

Summary
oo

Simple example with "∗"

▶ Does string 'bd' match regex 'ba∗d'?
  ▶ 'b' matches 'b'
  ▶ '' matches 'a∗'
  ▶ 'd' matches 'd'

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
o●oooo

Extended
oooooooo

Summary
oo

# Simple example with "*"

▶ Does string 'bd' match regex 'ba*d'?
  ▶ 'b' matches 'b'
  ▶ '' matches 'a*'
  ▶ 'd' matches 'd'
  Yes.

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

Repetition
0●0000

Extended
00000000

Summary
00

# Simple example with "*"

▶ Does string 'bd' match regex 'ba*d'?
  ▶ 'b' matches 'b'
  ▶ '' matches 'a*'
  ▶ 'd' matches 'd'

  Yes.
▶ Does string 'baaad' match regex 'ba*d'?

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

**Repetition**
0●0000

Extended
00000000

Summary
00

# Simple example with "∗"

- ▶ Does string 'bd' match regex 'ba∗d'?
  - ▶ 'b' matches 'b'
  - ▶ '' matches 'a∗'
  - ▶ 'd' matches 'd'

  Yes.
- ▶ Does string 'baaad' match regex 'ba∗d'?
  - ▶ 'b' matches 'b'
  - ▶ 'aaa' matches 'a∗'
  - ▶ 'd' matches 'd'

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

**Repetition**
0●0000

Extended
00000000

Summary
00

# Simple example with "∗"

▶ Does string 'bd' match regex 'ba∗d'?
   ▶ 'b' matches 'b'
   ▶ '' matches 'a∗'
   ▶ 'd' matches 'd'

Yes.

▶ Does string 'baaad' match regex 'ba∗d'?
   ▶ 'b' matches 'b'
   ▶ 'aaa' matches 'a∗'
   ▶ 'd' matches 'd'

Yes.

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

# Trickier example with "∗"

▶ Does string '7' match regex '[0-9]*[02468]'?

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

Repetition
000●000

Extended
00000000

Summary
00

# Trickier example with "∗"

- ▶ Does string '7' match regex '[0-9]∗[02468]'?
    - ▶ '7' matches '[0-9]∗' but '' does not match '[02468]'
    - ▶ '' matches '[0-9]∗' but '7' does not match '[02468]'

Introduction
ooo
Simple patterns
ooooo
Classes
oooooooo
Special
ooooooo
Repetition
oo●ooo
Extended
ooooooooo
Summary
oo

# Trickier example with "*"

- ▶ Does string '7' match regex '[0-9]*[02468]'?
    - ▶ '7' matches '[0-9]*' but '' does not match '[02468]'
    - ▶ '' matches '[0-9]*' but '7' does not match '[02468]'

    No.

Introduction
ooo

Simple patterns
ooooo

Classes
ooooooooo

Special
ooooooo

Repetition
oo●oooo

Extended
ooooooooo

Summary
oo

## Trickier example with "*"

- ▶ Does string '7' match regex '[0-9]*[02468]'?
  - ▶ '7' matches '[0-9]*' but '' does not match '[02468]'
  - ▶ '' matches '[0-9]*' but '7' does not match '[02468]'

  No.
- ▶ Does string '8' match regex '[0-9]*[02468]'?

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

Repetition
000●000

Extended
00000000

Summary
00

# Trickier example with "*"

- Does string '7' match regex '[0-9]*[02468]'?
    - '7' matches '[0-9]*' but '' does not match '[02468]'
    - '' matches '[0-9]*' but '7' does not match '[02468]'

    No.
- Does string '8' match regex '[0-9]*[02468]'?
    - '8' matches '[0-9]*' but '' does not match '[02468]'
    - '' matches '[0-9]*' and '8' matches '[02468]'

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○○

Special
○○○○○○○

**Repetition**
○○●○○○

Extended
○○○○○○○○○

Summary
○○

# Trickier example with "*"

- ▶ Does string '7' match regex '[0-9]*[02468]'?
  - ▶ '7' matches '[0-9]*' but '' does not match '[02468]'
  - ▶ '' matches '[0-9]*' but '7' does not match '[02468]'

  No.

- ▶ Does string '8' match regex '[0-9]*[02468]'?
  - ▶ '8' matches '[0-9]*' but '' does not match '[02468]'
  - ▶ '' matches '[0-9]*' and '8' matches '[02468]'

  Yes.

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

Repetition
00●000

Extended
00000000

Summary
00

# Trickier example with "*"

- ▶ Does string '7' match regex '[0-9]*[02468]'?
  - ▶ '7' matches '[0-9]*' but '' does not match '[02468]'
  - ▶ '' matches '[0-9]*' but '7' does not match '[02468]'

  No.

- ▶ Does string '8' match regex '[0-9]*[02468]'?
  - ▶ '8' matches '[0-9]*' but '' does not match '[02468]'
  - ▶ '' matches '[0-9]*' and '8' matches '[02468]'

  Yes.

- ▶ Does string '84' match regex '[0-9]*[02468]'?
  - ▶ '8' matches '[0-9]*' and '4' matches '[02468]'
  - ▶ '' matches '[0-9]*' and '8' matches '[02468]'
    but then we have an extra '4' that matches nothing

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

Repetition
000000

Extended
00000000

Summary
00

# Trickier example with "∗"

- ▶ Does string '7' match regex '`[0-9]*[02468]`'?
  - ▶ '7' matches '`[0-9]*`' but '' does not match '`[02468]`'
  - ▶ '' matches '`[0-9]*`' but '7' does not match '`[02468]`'

  No.

- ▶ Does string '8' match regex '`[0-9]*[02468]`'?
  - ▶ '8' matches '`[0-9]*`' but '' does not match '`[02468]`'
  - ▶ '' matches '`[0-9]*`' and '8' matches '`[02468]`'

  Yes.

- ▶ Does string '84' match regex '`[0-9]*[02468]`'?
  - ▶ '8' matches '`[0-9]*`' and '4' matches '`[02468]`'
  - ▶ '' matches '`[0-9]*`' and '8' matches '`[02468]`'
    but then we have an extra '4' that matches nothing

  Yes.

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
oo●ooo

Extended
oooooooo

Summary
oo

# Trickier example with "∗"

▶ Does string '7' match regex '[0-9]∗[02468]'?
  ▶ '7' matches '[0-9]∗' but '' does not match '[02468]'
  ▶ '' matches '[0-9]∗' but '7' does not match '[02468]'

  No.

▶ Does string '8' match regex '[0-9]∗[02468]'?
  ▶ '8' matches '[0-9]∗' but '' does not match '[02468]'
  ▶ '' matches '[0-9]∗' and '8' matches '[02468]'

  Yes.

▶ Does string '84' match regex '[0-9]∗[02468]'?
  ▶ '8' matches '[0-9]∗' and '4' matches '[02468]'
  ▶ '' matches '[0-9]∗' and '8' matches '[02468]'
    but then we have an extra '4' that matches nothing

  Yes.

  Remember: regular expressions match as much as they can

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○○

Special
○○○○○○○

**Repetition**
○○○●○○

Extended
○○○○○○○○

Summary
○○

# Some example regular expressions

▶ What does regular expression 'hello.*' mean?

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

**Repetition**
000●00

Extended
00000000

Summary
00

## Some example regular expressions

▶ What does regular expression 'hello.*' mean?
  ▶ 'hello' followed by any other characters

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

**Repetition**
ooo●oo

Extended
ooooooooo

Summary
oo

# Some example regular expressions

▶ What does regular expression 'hello.*' mean?
  ▶ 'hello' followed by any other characters
▶ Legal variable names in bash (and C, and Java...):

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

Repetition
000●00

Extended
00000000

Summary
00

## Some example regular expressions

- ▶ What does regular expression 'hello.*' mean?
  - ▶ 'hello' followed by any other characters
- ▶ Legal variable names in bash (and C, and Java...):

  ```
  [a-zA-Z_][a-zA-Z_0-9]*
  ```

## Some example regular expressions

- ▶ What does regular expression 'hello.*' mean?
  - ▶ 'hello' followed by any other characters
- ▶ Legal variable names in bash (and C, and Java...):

      [a-zA-Z_][a-zA-Z_0-9]*

- ▶ A string surrounded by double quotes
  (without allowing '\"' inside):

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○○

Special
○○○○○○○

**Repetition**
○○○●○○

Extended
○○○○○○○○

Summary
○○

## Some example regular expressions

- ▶ What does regular expression 'hello.*' mean?
  - ▶ 'hello' followed by any other characters
- ▶ Legal variable names in bash (and C, and Java...):

     [a-zA-Z_][a-zA-Z_0-9]*

- ▶ A string surrounded by double quotes
  (without allowing '\"' inside):

     "[^"]*"

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

**Repetition**
000●00

Extended
00000000

Summary
00

## Some example regular expressions

- ▶ What does regular expression 'hello.*' mean?
  - ▶ 'hello' followed by any other characters
- ▶ Legal variable names in bash (and C, and Java...):

    [a-zA-Z_][a-zA-Z_0-9]*

- ▶ A string surrounded by double quotes
  (without allowing '\"' inside):

    "[^"]*"

- ▶ Why not ".*" for strings?

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

Repetition
000●00

Extended
00000000

Summary
00

## Some example regular expressions

- ▶ What does regular expression 'hello.*' mean?
  - ▶ 'hello' followed by any other characters
- ▶ Legal variable names in bash (and C, and Java...):

    [a-zA-Z_][a-zA-Z_0-9]*

- ▶ A string surrounded by double quotes
  (without allowing '\"' inside):

    "[^"]*"

- ▶ Why not ".*" for strings?
  - ▶ Because "foo"bar" matches ".*"

# Specified number of repeats

> "$\backslash\{\ldots\backslash\}$": repeat the previous thing, some number of times
>
> $\backslash\{n\backslash\}$ : repeat exactly *n* times
>
> $\backslash\{n, m\backslash\}$ : repeat at least *n* and no more than *m* times
>
> $\backslash\{n,\backslash\}$ : repeat at least *n* times
>
> Not supported everywhere

Note: "$*$" is the same as "$\backslash\{0,\backslash\}$"

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
oooooo●

Extended
oooooooo

Summary
oo

# Some examples

▶ Date strings (`mm/dd/yyyy` or `mm-dd-yyyy`)

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

**Repetition**
00000●

Extended
00000000

Summary
00

## Some examples

▶ Date strings (mm/dd/yyyy or mm-dd-yyyy)

[0-9]\{2\}[/-][0-9]\{2\}[/-][0-9]\{4\}

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

**Repetition**
00000●

Extended
00000000

Summary
00

## Some examples

▶ Date strings (mm/dd/yyyy or mm-dd-yyyy)

   [0-9]\{2\}[/-][0-9]\{2\}[/-][0-9]\{4\}

   Note that mm/dd-yyyy and mm-dd/yyyy will match also

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○○

Special
○○○○○○○

**Repetition**
○○○○○●

Extended
○○○○○○○○

Summary
○○

# Some examples

▶ Date strings (`mm/dd/yyyy` or `mm-dd-yyyy`)

$\quad$ `[0-9]\{2\}[/-][0-9]\{2\}[/-][0-9]\{4\}`

$\quad$ Note that `mm/dd-yyyy` and `mm-dd/yyyy` will match also

▶ City, state, zip code

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○○

Special
○○○○○○○

**Repetition**
○○○○○●

Extended
○○○○○○○○

Summary
○○

## Some examples

▶ Date strings (mm/dd/yyyy or mm-dd-yyyy)

    [0-9]\{2\}[/-][0-9]\{2\}[/-][0-9]\{4\}

    Note that  mm/dd-yyyy  and  mm-dd/yyyy  will match also

▶ City, state, zip code

    [- A-Za-z]*, [A-Z]\{2\}, [0-9]\{5\}

# Extended regular expressions

- ▶ Have additional features and a few small differences
- ▶ Used by egrep or grep -e
- ▶ Used by sed and awk

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

Repetition
000000

Extended
0●000000

Summary
00

# Simple changes in extended regular expressions

- ▶ It is not necessary to escape the braces. E.g.:
    - ▶ Use '$\{n\}$' instead of '$\backslash\{n\backslash\}$'
- ▶ '+' is shorthand for '$\{1,\}$' (repeat one or more times)
- ▶ '?' is shorthand for '$\{0,1\}$' (previous thing is optional)

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
oooooo

Extended
ooooooo

Summary
oo

# Examples using egrep

```
prompt$
```

# Examples using egrep

```
prompt$ echo 'a' | egrep '^ab*$'
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
oooooo

Extended
oo●ooooo

Summary
oo

# Examples using egrep

```
prompt$ echo 'a' | egrep '^ab*$'
a
prompt$
```

# Examples using egrep

```
prompt$ echo 'a' | egrep '^ab*$'
a
prompt$ echo 'a' | egrep '^ab+$'
```

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○○

Special
○○○○○○○

Repetition
○○○○○○

**Extended**
○○●○○○○○

Summary
○○

# Examples using egrep

```
prompt$ echo 'a' | egrep '^ab*$'
a
prompt$ echo 'a' | egrep '^ab+$'
prompt$
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
oooooo

Extended
ooo●oooo

Summary
oo

# Examples using egrep

```
prompt$ echo 'a' | egrep '^ab*$'
a
prompt$ echo 'a' | egrep '^ab+$'
prompt$ echo 'ab' | egrep '^ab+$'
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
oooooo

Extended
oo●oooooo

Summary
oo

# Examples using egrep

```
prompt$ echo 'a' | egrep '^ab*$'
a
prompt$ echo 'a' | egrep '^ab+$'
prompt$ echo 'ab' | egrep '^ab+$'
ab
prompt$ █
```

# Examples using `egrep`

```
prompt$ echo 'a' | egrep '^ab*$'
a
prompt$ echo 'a' | egrep '^ab+$'
prompt$ echo 'ab' | egrep '^ab+$'
ab
prompt$ echo 'abbb' | egrep '^ab+$'
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
oooooo

Extended
ooooooo

Summary
oo

# Examples using `egrep`

```
prompt$ echo 'a' | egrep '^ab*$'
a
prompt$ echo 'a' | egrep '^ab+$'
prompt$ echo 'ab' | egrep '^ab+$'
ab
prompt$ echo 'abbb' | egrep '^ab+$'
abbb
prompt$ 
```

Introduction
ooo

Simple patterns
ooooo

Classes
ooooooo

Special
ooooooo

Repetition
oooooo

Extended
ooooooo

Summary
oo

# Examples using `egrep`

```
prompt$ echo 'a' | egrep '^ab*$'
a
prompt$ echo 'a' | egrep '^ab+$'
prompt$ echo 'ab' | egrep '^ab+$'
ab
prompt$ echo 'abbb' | egrep '^ab+$'
abbb
prompt$ echo 'abbb' | egrep '^ab?$'
```

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
oooooo

Extended
ooooooo

Summary
oo

# Examples using `egrep`

```
prompt$ echo 'a' | egrep '^ab*$'
a
prompt$ echo 'a' | egrep '^ab+$'
prompt$ echo 'ab' | egrep '^ab+$'
ab
prompt$ echo 'abbb' | egrep '^ab+$'
abbb
prompt$ echo 'abbb' | egrep '^ab?$'
prompt$
```

Introduction
ooo

Simple patterns
ooooo

Classes
ooooooo

Special
ooooooo

Repetition
oooooo

Extended
ooooooo

Summary
oo

# Examples using `egrep`

```
prompt$ echo 'a' | egrep '^ab*$'
a
prompt$ echo 'a' | egrep '^ab+$'
prompt$ echo 'ab' | egrep '^ab+$'
ab
prompt$ echo 'abbb' | egrep '^ab+$'
abbb
prompt$ echo 'abbb' | egrep '^ab?$'
prompt$ echo 'ab' | egrep '^ab?$'
```

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

Repetition
000000

Extended
00●00000

Summary
00

# Examples using `egrep`

```
prompt$ echo 'a' | egrep 'ˆab*$'
a
prompt$ echo 'a' | egrep 'ˆab+$'
prompt$ echo 'ab' | egrep 'ˆab+$'
ab
prompt$ echo 'abbb' | egrep 'ˆab+$'
abbb
prompt$ echo 'abbb' | egrep 'ˆab?$'
prompt$ echo 'ab' | egrep 'ˆab?$'
ab
prompt$
```

Introduction
ooo

Simple patterns
ooooo

Classes
ooooooooo

Special
ooooooo

Repetition
oooooo

Extended
oo●ooooo

Summary
oo

# Examples using `egrep`

```
prompt$ echo 'a' | egrep '^ab*$'
a
prompt$ echo 'a' | egrep '^ab+$'
prompt$ echo 'ab' | egrep '^ab+$'
ab
prompt$ echo 'abbb' | egrep '^ab+$'
abbb
prompt$ echo 'abbb' | egrep '^ab?$'
prompt$ echo 'ab' | egrep '^ab?$'
ab
prompt$ echo 'a' | egrep '^ab?$'
```

# Examples using `egrep`

```
prompt$ echo 'a' | egrep '^ab*$'
a
prompt$ echo 'a' | egrep '^ab+$'
prompt$ echo 'ab' | egrep '^ab+$'
ab
prompt$ echo 'abbb' | egrep '^ab+$'
abbb
prompt$ echo 'abbb' | egrep '^ab?$'
prompt$ echo 'ab' | egrep '^ab?$'
ab
prompt$ echo 'a' | egrep '^ab?$'
a
prompt$ █
```

# Significant extensions

- ▶ There are 2 more things to discuss
- ▶ These give regular expressions much more power
- ▶ They also make things much more complicated

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
oooooo

Extended
oooo●oooo

Summary
oo

# Significant extensions

- ▶ There are 2 more things to discuss
- ▶ These give regular expressions much more power
- ▶ They also make things much more complicated

- ▶ Now is a good time to ask questions before things get crazy (or review if you are reading this at home)

# Grouping

Regular expressions may be grouped using ()

# Grouping

Regular expressions may be grouped using ()

▶ The "previous thing" for a repeat operator may be a regex

# Grouping

Regular expressions may be grouped using ()

▶ The "previous thing" for a repeat operator may be a regex
▶ We can repeat entire subexpressions
▶ We can make entire subexpressions optional

Simple example:

    ab+ means:

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
oooooo

**Extended**
ooooo●ooo

Summary
oo

# Grouping

### Regular expressions may be grouped using ()

► The "previous thing" for a repeat operator may be a regex
► We can repeat entire subexpressions
► We can make entire subexpressions optional

Simple example:

ab+ means: 'a' followed by one or more 'b's

► Strings 'ab', 'abb', 'abbb', ..., match this

(ab)+ means:

# Grouping

## Regular expressions may be grouped using ()

▶ The "previous thing" for a repeat operator may be a regex

▶ We can repeat entire subexpressions

▶ We can make entire subexpressions optional

Simple example:

ab+ means: 'a' followed by one or more 'b's

▶ Strings 'ab', 'abb', 'abbb', ..., match this

(ab)+ means: a sequence of one or more 'ab's

▶ Strings 'ab', 'abab', 'ababab', ..., match this

(ab+)+ means:

# Grouping

## Regular expressions may be grouped using ()

▶ The "previous thing" for a repeat operator may be a regex
▶ We can repeat entire subexpressions
▶ We can make entire subexpressions optional

Simple example:

ab+ means: 'a' followed by one or more 'b's

▶ Strings 'ab', 'abb', 'abbb', ..., match this

(ab)+ means: a sequence of one or more 'ab's

▶ Strings 'ab', 'abab', 'ababab', ..., match this

(ab+)+ means:

▶ One or more strings matching ab+, concatenated

# Grouping

## Regular expressions may be grouped using ()

► The "previous thing" for a repeat operator may be a regex

► We can repeat entire subexpressions

► We can make entire subexpressions optional

Simple example:

ab+ means: 'a' followed by one or more 'b's

    ► Strings 'ab', 'abb', 'abbb', ..., match this

(ab)+ means: a sequence of one or more 'ab's

    ► Strings 'ab', 'abab', 'ababab', ..., match this

(ab+)+ means:

    ► One or more strings matching ab+, concatenated

    ► Start with 'a', end with 'b', never 2 'a's together

# Serious grouping examples

- ▶ Zip code with optional "extra 4":

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○○

Special
○○○○○○○

Repetition
○○○○○○

**Extended**
○○○○○●○○

Summary
○○

## Serious grouping examples

▶ Zip code with optional "extra 4":

[0-9]{5}(-[0-9]{4})?

## Serious grouping examples

▶ Zip code with optional "extra 4":

  `[0-9]{5}(-[0-9]{4})?`

▶ Well–formed strings, where '\"' is allowed inside

## Serious grouping examples

▶ Zip code with optional "extra 4":

  `[0-9]{5}(-[0-9]{4})?`

▶ Well–formed strings, where '\"' is allowed inside

  `"([^"]?(\\")?)*"`

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
oooooo

Extended
oooooo●oo

Summary
oo

# Serious grouping examples

▶ Zip code with optional "extra 4":

    [0-9]{5}(-[0-9]{4})?

▶ Well–formed strings, where '\"' is allowed inside

    "([^"]?(\\")?)*"

▶ A "real value" constant as allowed in C or Java

Introduction
ooo
Simple patterns
ooooo
Classes
oooooooo
Special
ooooooo
Repetition
oooooo
Extended
ooooo●oo
Summary
oo

# Serious grouping examples

▶ Zip code with optional "extra 4":

    [0-9]{5}(-[0-9]{4})?

▶ Well–formed strings, where '\"' is allowed inside

    "([^"]?(\\")?)*"

▶ A "real value" constant as allowed in C or Java

    -?[0-9]+(\.[0-9]+)?([eE]-?[0-9]+)?

Introduction
ooo

Simple patterns
ooooo

Classes
oooooooo

Special
ooooooo

Repetition
oooooo

**Extended**
ooooooo●o

Summary
oo

# Choosing between expressions

## '|' means "or"

- ▶ Not needed for single characters — use [] instead
- ▶ Makes sense when we group things

Simple example:

(ab)|(cd)|(ef) means: 'ab' or 'cd' or 'ef'

```
prompt$
```

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○○

Special
○○○○○○○

Repetition
○○○○○○

**Extended**
○○○○○○●○

Summary
○○

# Choosing between expressions

### '|' means "or"

- ▶ Not needed for single characters — use [] instead
- ▶ Makes sense when we group things

Simple example:

(ab)|(cd)|(ef) means: 'ab' or 'cd' or 'ef'

```
prompt$ echo 'cf' | egrep '^((ab)|(cd)|(ef))$'
```

Introduction
ooo
Simple patterns
ooooo
Classes
ooooooooo
Special
ooooooo
Repetition
oooooo
**Extended**
ooooooo●o
Summary
oo

# Choosing between expressions

### '|' means "or"

- ► Not needed for single characters — use [] instead
- ► Makes sense when we group things

Simple example:

(ab)|(cd)|(ef) means: 'ab' or 'cd' or 'ef'

```
prompt$ echo 'cf' | egrep '^((ab)|(cd)|(ef))$'
prompt$ █
```

Introduction
○○○

Simple patterns
○○○○○

Classes
○○○○○○○○

Special
○○○○○○○

Repetition
○○○○○○

Extended
○○○○○○●○

Summary
○○

# Choosing between expressions

## '|' means "or"

▶ Not needed for single characters — use [] instead
▶ Makes sense when we group things

Simple example:

(ab)|(cd)|(ef) means: 'ab' or 'cd' or 'ef'

```
prompt$ echo 'cf' | egrep '^((ab)|(cd)|(ef))$'
prompt$ echo 'ab' | egrep '^((ab)|(cd)|(ef))$'
```

# Choosing between expressions

## '|' means "or"

- ▶ Not needed for single characters — use [] instead
- ▶ Makes sense when we group things

Simple example:

(ab)|(cd)|(ef) means: 'ab' or 'cd' or 'ef'

```
prompt$ echo 'cf' | egrep '^((ab)|(cd)|(ef))$'
prompt$ echo 'ab' | egrep '^((ab)|(cd)|(ef))$'
ab
prompt$ █
```

Introduction
ooo
Simple patterns
ooooo
Classes
oooooooo
Special
ooooooo
Repetition
oooooo
Extended
ooooooo●
Summary
oo

# Final examples

▶ Date strings (mm/dd/yyyy or mm-dd-yyyy)

# Final examples

▶ Date strings (mm/dd/yyyy or mm-dd-yyyy)

  [0-9]{2}(-[0-9]{2}-)|(/[0-9]{2}/)[0-9]{4}

▶ Date strings where:
  ▶ The month is one or two digits, and is a legal month
  ▶ The day is one or two digits
  ▶ The year is two or four digits
  ▶ '/' separators only (let's make it easy)

## Final examples

▶ Date strings (`mm/dd/yyyy` or `mm-dd-yyyy`)

`[0-9]{2}(-[0-9]{2}-)|(/[0-9]{2}/)[0-9]{4}`

▶ Date strings where:
  ▶ The month is one or two digits, and is a legal month
  ▶ The day is one or two digits
  ▶ The year is two or four digits
  ▶ '/' separators only (let's make it easy)

`((1[0-2])|[1-9])/[0-9]{1,2}/([0-9]{2}){1,2}`

Introduction
000

Simple patterns
00000

Classes
00000000

Special
0000000

Repetition
000000

Extended
00000000

Summary
●○

## Basic regular expressions

ordinary character  match itself

[] match one character from the list

- set a range (unless it is first or last)
- ^ invert the list (if it appears first)

. match any character

\ convert meta to ordinary character

^ imaginary beginning of line character

$ imaginary end of line character

* match previous character zero or more times

\{n, m\} match previous character between *n* and *m* times

Introduction
ooo

Simple patterns
ooooo

Classes
ooooooooo

Special
ooooooo

Repetition
oooooo

Extended
ooooooooo

Summary
o●

# Extended regular expressions

Same as "basic" except change or add:

|             |                                                    |
|------------:|----------------------------------------------------|
| ( )         | group patterns                                     |
| *           | match previous pattern zero or more times          |
| +           | match previous pattern one or more times           |
| ?           | match previous pattern zero or one times           |
| $\{n,m\}$   | match previous pattern between $n$ and $m$ times   |
| \|          | choose between patterns                            |

Introduction
ooo
Simple patterns
ooooo
Classes
oooooooo
Special
ooooooo
Repetition
oooooo
Extended
oooooooo
Summary
oo

End of lecture