

Web Scraping

Extracting data from the web

APIs and beyond



Scott Chamberlain
Karthik Ram
Garrett Grolemund

June 2016

What if data is on a web page but there is no API?

www.imdb.com/title/tt2294629/

The screenshot shows a web browser window displaying the IMDb movie page for "Frozen (2013)". The page includes the movie's title, rating (7.6/10), and a play button for a video thumbnail. A prominent red Toyota Prius advertisement is overlaid on the page, featuring the text "THE ALL-NEW PRIUS EXHILARATING DRIVING DYNAMICS. TIME FOR HYBRIDS TO HAVE FUN." and a "LEARN MORE" button. The browser interface at the top shows the URL "www.imdb.com/title/tt2294629/" and various browser extensions.

IMDb Picks: June

Strategy

Garrett

IMDb Frozen (2013) - IMDb

www.imdb.com/title/tt2294629/

Apps SelectorGadget Other Bookmarks

People who liked this also liked...

Tangled (2010)
PG Animation | Adventure | Comedy
7.8/10

The magically long-haired Rapunzel has spent her entire life in a tower, but now that a runaway thief has stumbled upon her, she is about to discover the world for the first time, and who she really is.

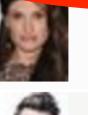
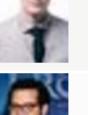
Add to Watchlist

Next >

◀ Prev 6 Next 6 ▶

Cast

Cast overview, first billed only:

	Kristen Bell	... Anna (voice)
	Idina Menzel	... Elsa (voice)
	Jonathan Groff	... Kristoff (voice)
	Josh Gad	... Olaf (voice)
	Santino Fontana	... Hans (voice)
	Alan Tudyk	... Duke (voice)
	Ciarán Hinds	... Pabbie / Grandpa (voice)

Garrett

IMDb Frozen (2013) - IMDb

view-source:www.imdb.com

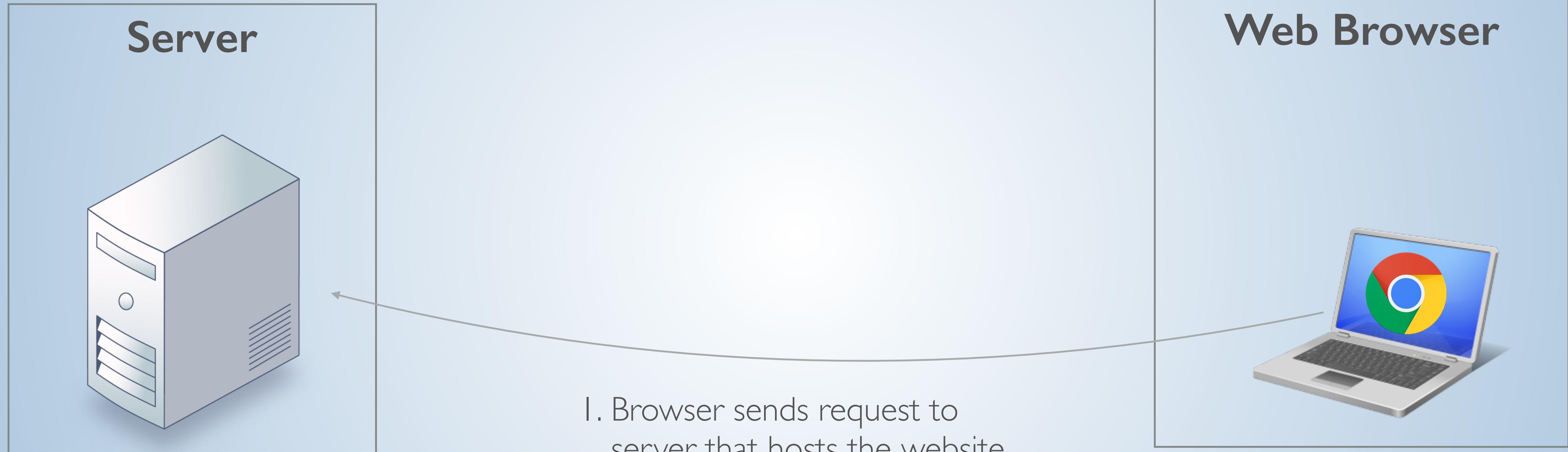
view-source:www.imdb.com/title/tt2294629/

Apps SelectorGadget Other Bookmarks

```
3933 <table class="cast_list">
3934 <tr><td colspan="4" class="castlist_label">Cast
3935 overview, first billed only:</td></tr>
3936 <tr class="odd">
3937 <td class="primary_photo">
3938 <a href="/name/nm0068338/?ref_=tt_cl_i1"
3939 ></a> </td>
3940 <td class="itemprop" itemprop="actor"
3941 itemscope itemtype="http://schema.org/Person">
3942 <a href="/name/nm0068338/?ref_=tt_cl_t1"
3943 itemprop='url'><span class="itemprop"
3944 itemprop="name">Kristen Bell</span>
3945 </a> </td>
3946 <td class="ellipsis">
3947 ...
3948 </td>
3949 <td class="character">
3950 <div>
3951 <a href="/character/ch0307445/?ref_=tt_cl_t1" >Anna</a>
3952 (voice)
3953
```

A large grey arrow points from the left browser window to the right one.

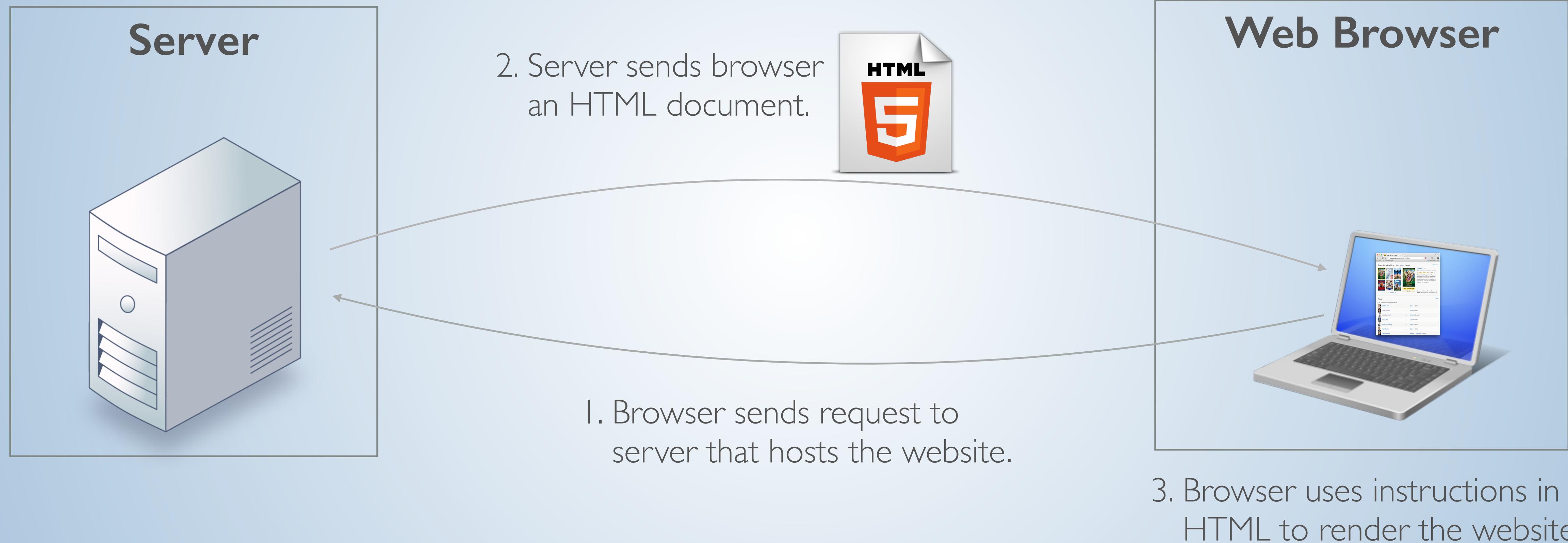
HTML (Review)



HTML (Review)



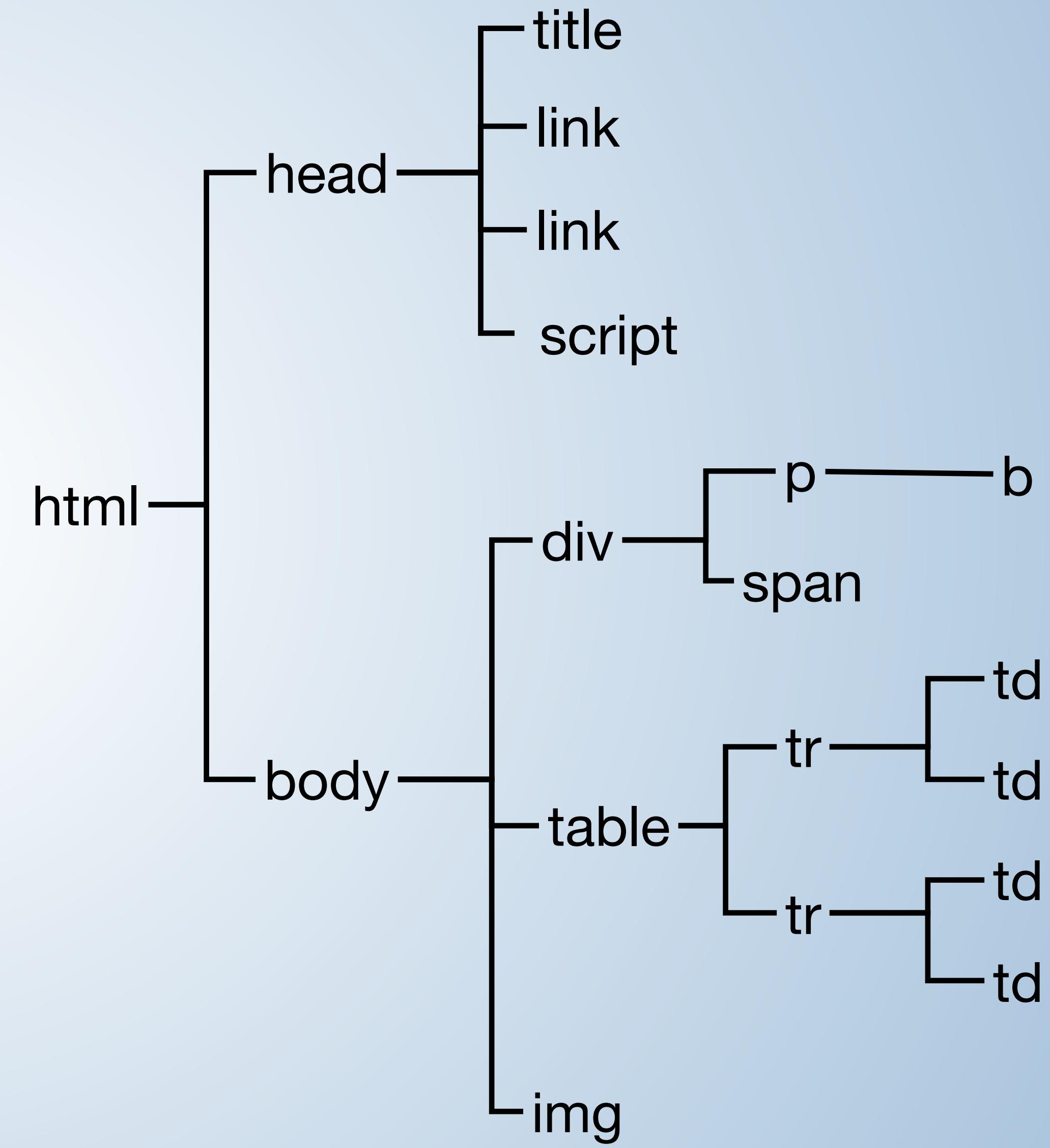
HTML (Review)



HTML (Review)



```
<html>
  <head>
    <title>Title</title>
    <link rel="icon" type="icon" href="http://a" />
    <link rel="icon" type="icon" href="http://b" />
    <script type="text/javascript">
      var ue_t0=window.ue_t0||+new Date();
    </script>
  </head>
  <body>
    <div>
      <p>Click <b>here</b> now.</p>
      <span>Frozen</span>
    </div>
    <table style="width:100%">
      <tr>
        <td>Kristen</td>
        <td>Bell</td>
      </tr>
      <tr>
        <td>Idina</td>
        <td>Menzel</td>
      </tr>
    </table>
    
  </body>
</html>
```



HTML (Review)

Each element in the page is created by a tag.

```
<a href="http://github.com">GitHub</a>
```

tag name

attribute
(name)

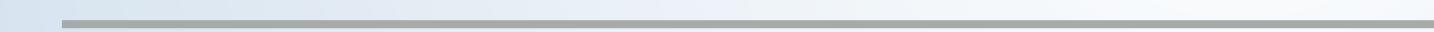
attribute
(value)

content

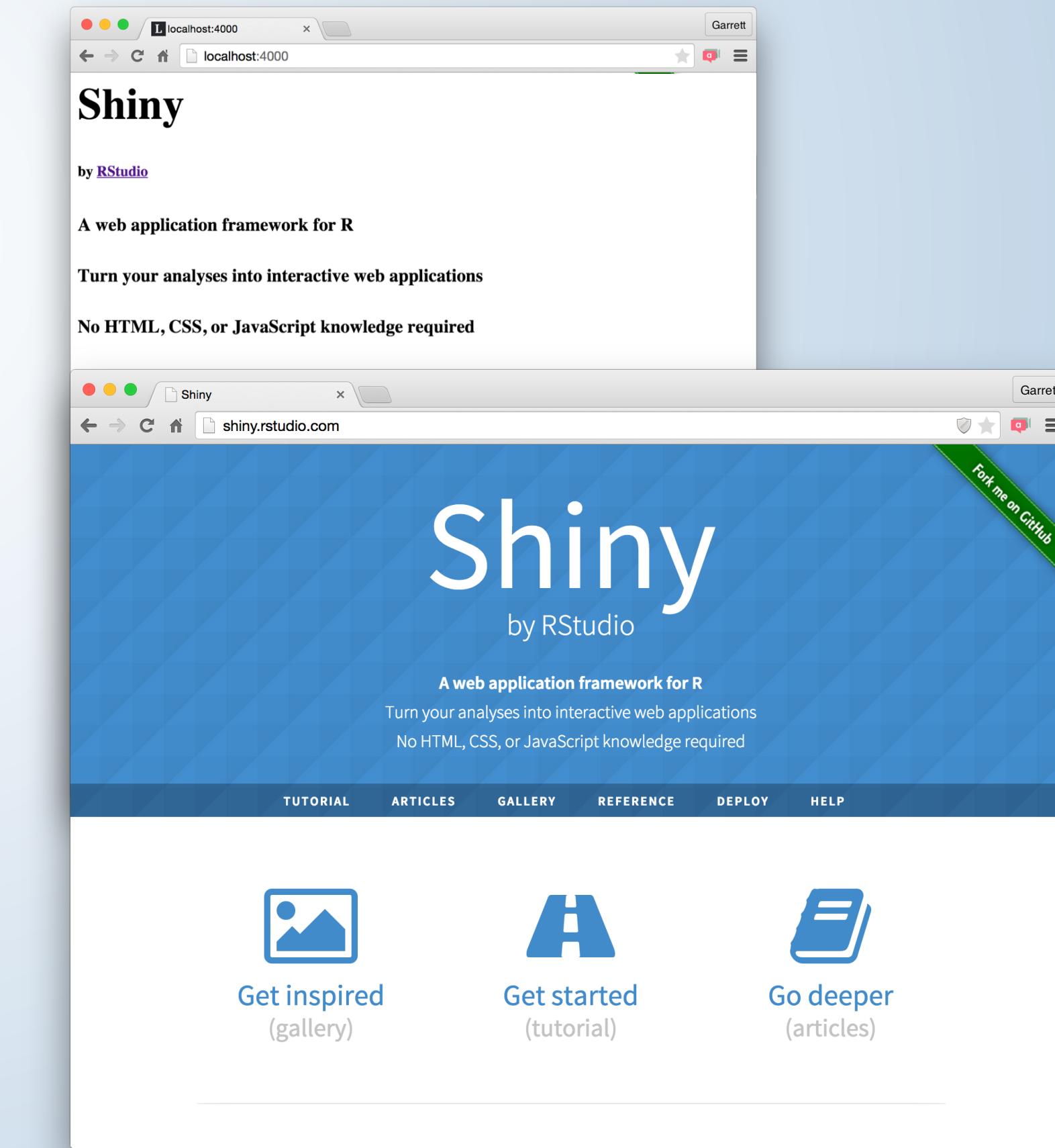
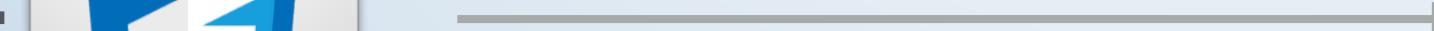
css selectors

CSS (*Review*)

Cascading Style Sheets (CSS) are a framework for customizing the appearance of elements in a web page.



+



CSS (*Review*)



localhost:4000 Garrett

Shiny

by RStudio

A web application framework for R

Turn your analyses into interactive web applications

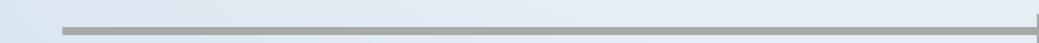
No HTML, CSS, or JavaScript knowledge required

- [Tutorial](#)
- [Articles](#)
- [Gallery](#)
- [Reference](#)
- [Deploy](#)
- [Help](#)

[Get inspired](#)
(gallery)

[Get started](#)
(tutorial)

[Go deeper](#)
(articles)



localhost:4000 Garrett

Shiny

by RStudio

A web application framework for R

Turn your analyses into interactive web applications

No HTML, CSS, or JavaScript knowledge required

[TUTORIAL](#) [ARTICLES](#) [GALLERY](#) [REFERENCE](#) [DEPLOY](#) [HELP](#)

 [Get inspired](#)
(gallery)

 [Get started](#)
(tutorial)

 [Go deeper](#)
(articles)

Fork me on GitHub

CSS (Review)



```
span {  
    color: #ffffff;  
}  
  
.num {  
    color: #a8660d;  
}  
  
table.data {  
    width: auto;  
}  
  
#firstname {  
    background-color: yellow;  
}
```

← selector

← styling

← selector

← styling

← selector

← styling

← selector

← styling

CSS (*Review*)

A CSS script describes an element by its tag, class, and/or ID.

```
<span class="bigname" id="shiny">Shiny</span>
```

tag name

class
(optional)

id
(optional)

CSS (*Review*)

A CSS script describes an element by its tag, class, and/or ID.

```
<span class="bigname" id="shiny">Shiny</span>
```

```
span
```

CSS selector for **ALL** elements with:

- the **span tag**

CSS (*Review*)

A CSS script describes an element by its tag, class, and/or ID.

```
<span class="bigname" id="shiny">Shiny</span>
```

```
.bigname
```

CSS selector for **ALL** elements with:

- the **bigname class**

CSS (*Review*)

A CSS script describes an element by its tag, class, and/or ID.

```
<span class="bigname" id="shiny">Shiny</span>
```

```
span.bigname
```

CSS selector for **ALL** elements with:

- the **span tag**

AND

- the **bigname class**

CSS (*Review*)

A CSS script describes an element by its tag, class, and/or ID.

```
<span class="bigname" id="shiny">Shiny</span>
```

```
#shiny
```

CSS selector for **ALL** elements with:

- the **shiny** id

Recap

Extract information from the HTML document.

Identify information to extract with CSS selectors.

rvest

rvest



A package that makes it easy to extract
info from a webpage.

```
install.packages("rvest")
```

* This will also install `xml2`, a package that `rvest` relies on.

Basic Workflow

1. Download the HTML and turn it into an XML file with `read_html()`

```
library(rvest)  
frozen <- read_html("http://www.imdb.com/title/tt2294629/")
```

read_html

URL

* `read_html()` comes in the `xml2` package

To examine contents

frozen

html_structure(frozen)

as_list(frozen)

xml_children(frozen)

xml_children(frozen)[[2]]

xml_contents(xml_children(frozen)[[2]])

Basic Workflow

1. Download the HTML and turn it into an XML file with `read_html()`
2. Extract specific nodes with `html_nodes()`

```
itals <- html_nodes(frozen, "em")
```

XML

emphasized tag

CSS selector

Basic Workflow

1. Download the HTML and turn it into an XML file with `read_html()`
2. Extract specific nodes with `html_nodes()`
3. Extract content from nodes with `html_text()`,
`html_name()`, `html_attrs()`, `html_children()`,
`html_table()`

Recap

1. Download the HTML and turn it into an XML file with `read_html()`
2. Extract specific nodes with `html_nodes()`
3. Extract content from nodes with `html_text()`,
`html_name()`, `html_attrs()`, `html_children()`,
`html_table()`

selectorGadget

selectorGadget

A GUI tool to identify CSS selector combinations



To Install

1. Run `vignette("selectorgadget")`
2. Drag `Selectorgadget` link into your browser's bookmark bar

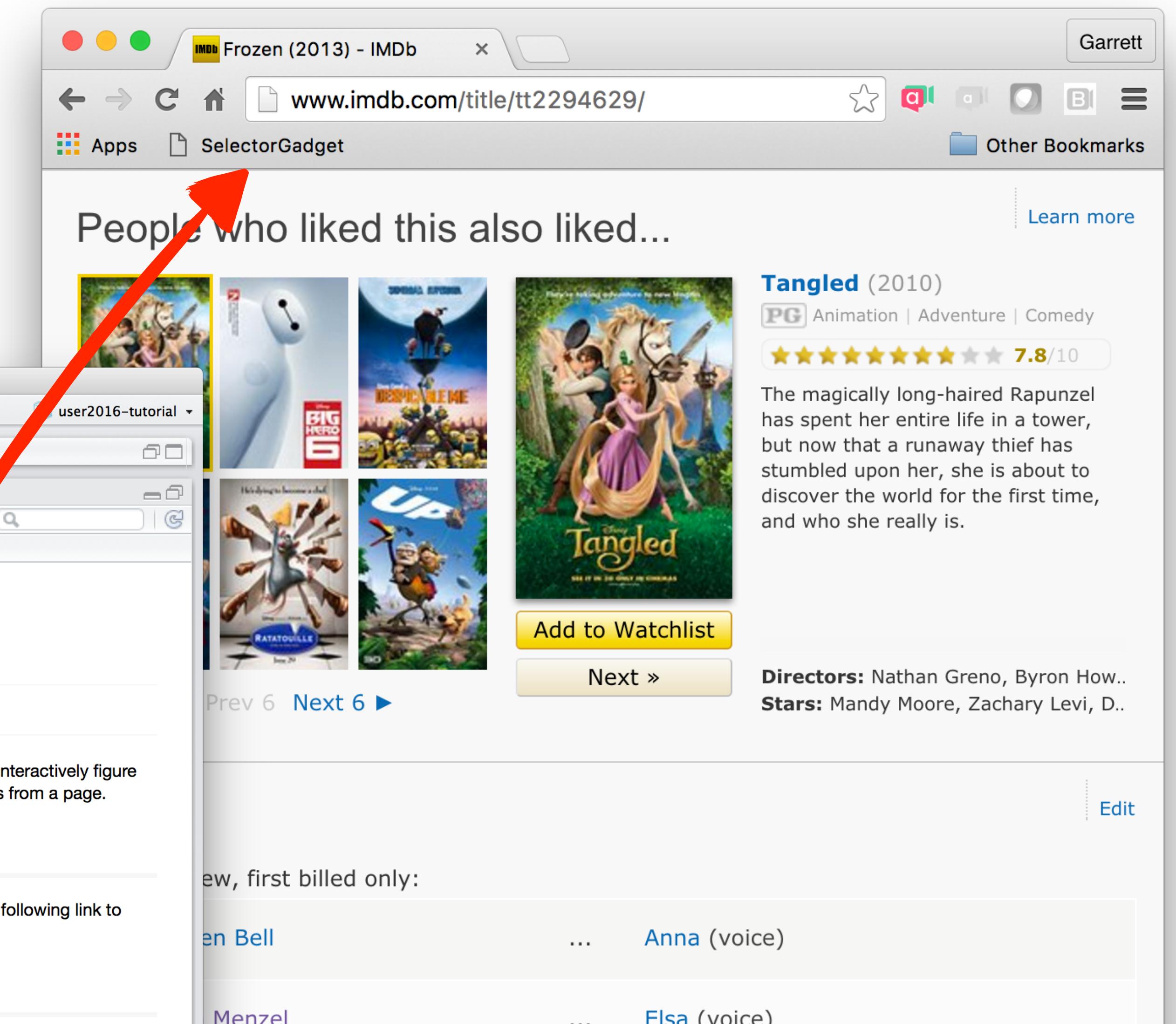
The screenshot shows the RStudio interface. In the top-left pane, there is an R script named 'scraping-outline.Rmd' with the following code:

```
1 # Read in web page
2
3 library(rvest)
4 frozen <- read_html("http://www.imdb.com/title/tt2294629")
5
6 # Look at web page
7
8 frozen
9 html_structure(frozen)
10 as_list(frozen)
11
12 xml_children(frozen)
13 xml_children(frozen)[[2]]
```

In the bottom-left pane, the console shows the command:

```
> vignette("selectorgadget")
```

The main pane displays the 'Selectorgadget' vignette page. A red circle highlights the instruction: "To install it, open this page in your browser and then drag the following link to your bookmark bar: [Selectorgadget](#)". A red arrow points from this instruction to the bookmark bar area of the browser window in the next image.



To Use

1. **Navigate** to a webpage
2. **Open** the SelectorGadget bookmark
3. **Click** on item to scrape
4. **Click** on yellow items You do not want to scrape
5. **Click** on additional items that you do want to scrape
6. **Copy** selector to use with `html_nodes()`

.fa-bolt

Clear (1) Toggle Position XPath Help X

CSS selector to use

start over

move gadget

show XPath

help

close gadget

tables

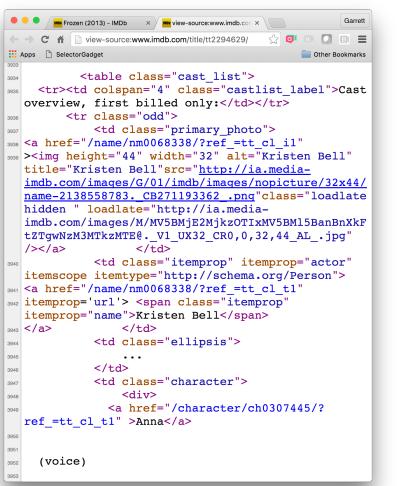
Tables

Use `html_table()` to scrape whole tables of data **as
a data frame.**

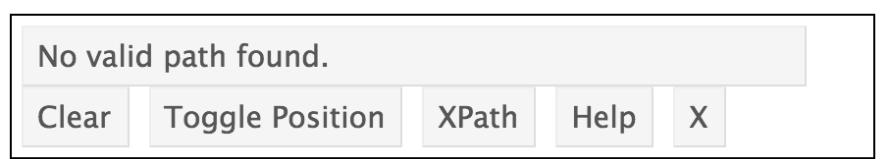
```
tables <- html_nodes(kw, css = "table")
html_table(tables, header = TRUE)[[2]]
```

Recap

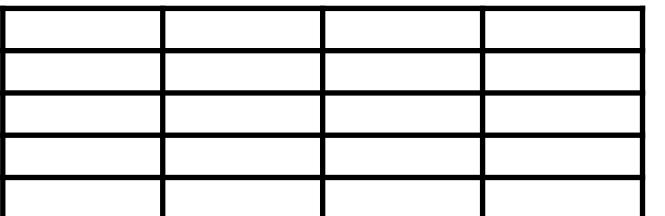
Pull from HTML, use HTML/CSS structure



read_html() **html_nodes()** **html_text()**,



selectorGadget for finding useful selector combinations



html_table() for tables

thank you