

Dr. Dobb's Journal

May 2013

Next

Supercharging ASP.NET Web Form Apps

ALSO INSIDE

[Microsoft TypeScript:
The Lay of the Land >>](#)

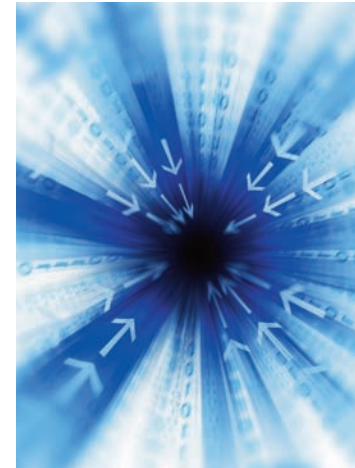
[Responsive Web Design >>](#)

[From the Vault: Building Scalable Web
Architecture and Distributed Systems >>](#)

Dr. Dobb's Journal

CONTENTS

May 2013

**COVER STORY****6 Supercharging ASP.NET Web Form Apps****By Charles Nurse**

The majority of ASP.NET websites are built using the ASP.NET Web Forms model. It provides a simple, stateful abstraction over the stateless HTTP protocol, and focuses on server-side processing. Recently, websites have focused more on a client-side model with the use of JavaScript, AJAX, and REST-style services. With the introduction of ASP.NET Web API in ASP.NET 4.5, it is now possible to add advanced client-side functionality to existing sites, while retaining many of the benefits of the Web Forms model.

EDITORIAL**3 Responsive Web Design: A Device-Oriented Salvation****By Dino Esposito**

Responsive Web Design is providing an excellent solution to sites that need to respond to mobile users, but don't want to tailor apps to an infinite number of devices and form factors.

FEATURES**12 Microsoft TypeScript: The Lay of the Land****By Gaston Hillar**

Redmond's new language adds types to JavaScript syntax, but it compiles to pure JavaScript. How much does the typing help and how much clutter does it add to the resulting JavaScript?

20 From the Vault: Building Scalable Web Architecture and Distributed Systems**By Kate Matsudaira**

What exactly does it mean to build and operate a scalable website or application? At a primitive level it's just connecting users with remote resources via the Internet — the part that makes it scalable is that the resources, or access to those resources, are distributed across multiple servers. Like most things in life, taking the time to plan ahead when building a Web service can help in the long run.

39 Links

Snapshots of the most interesting items on drdobbs.com including MongoDB with C# and NoSQL options compared.

More on DrDobbs.com**Hadoop: The Lay of the Land**

The core map-reduce framework for big data consists of several interlocking technologies. This first installment of our tutorial explains what Hadoop does and how the pieces fit together.

<http://www.drdobbs.com/240150854>

Pulling Back from Windows 8

The lack of hardware sales, the deep dissatisfaction among consumers, and the lack of interest in Win8-specific development suggest that Microsoft needs an immediate restart.

<http://www.drdobbs.com/240153041>

Optimizing a Program Means Making It Run Faster, Right?

People often use optimization to mean changing a program in ways that they think will make it run faster, but that casual definition is far from complete.

<http://www.drdobbs.com/240152663>

The Pragmatic Side of Using Big Data

A successful big data project is far less about correlating disparate data structures and far more about articulating a testable hypothesis, designing that test, and evaluating the results.

<http://www.drdobbs.com/240152350>

Roll Your Own Analog Input

Many systems (including systems on a chip) have methods to generate and read different analog quantities.

<http://www.drdobbs.com/240152547>

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

Responsive Web Design: A Device-Oriented Salvation

Responsive Web design is providing an excellent solution to sites that need to respond to mobile users, but don't want to tailor apps to an infinite number of devices and form factors.

By Dino Esposito

The vast majority of websites are neither optimized for nor adapted to mobile devices. Furthermore, a significant number of developers believe that making a website mobile-friendly is, overall, an easy task that just requires a smarter CSS file on top of the same old pages.

The major issue of mobile site development is the selection of use-cases that are most appropriate for the particular business scenario. Before mobile devices became so popular, developers built websites around a single set of use cases. At most, developers had to be concerned about how pages were rendered by different Web browsers. Ensuring that the same content was rendered the same way across all browsers has been a source of significant headaches for developers:

The ability to shield developers from the nuances of different browsers was one of the reasons for the rapid adoption of the jQuery library.

Mobile devices add a new dimension to the problem of identifying the right use-cases for a site: A mobile site may need to have its own use cases, which are inherently different from those of the full desktop site. In addition, the range of different devices available out there is large and can be measured in the order of thousands. Building a mobile site can be a walk in the park if you don't care about the requesting devices. However, it can get really problematic if you intend to serve each device some tailor-made content.

Web developers learned a key lesson over the years: Always use feature detection over browser detection. Web pages built around

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

the detected capabilities of the device do not need to be updated when a new device comes out. Conversely, Web pages specifically created for a particular browser (or device) are constantly subject to being extended and/or fixed when new device versions are released.

A relatively new approach to website development, Responsive Web Design (RWD) derives perfectly from the philosophy that detecting device features is smarter than detecting browsers.

Responsive Web Design

RWD is a design philosophy centered on making Web pages highly adaptive and capable of providing a good viewing experience regardless of the device details. The beauty of RWD is that it allows you to write a website once and view it effectively on a wide range of devices including laptops, smartphones, tablets, and more.

Technically speaking, RWD is based on two pillars: CSS media queries and proportional grids. CSS media queries let developers bind style sheets to conditions that the browser dynamically evaluates when some system event takes place. For example, an RWD-enabled site can reshape its content based on the size of the browser window. A CSS media query is an expression assigned to the media attribute of the `LINK` element:

```
<link type="text/css"
      rel="stylesheet"
      href="site.tablet.css"
      media="only screen and (max-width: 800px)">
```

The CSS file referenced in the code snippet is applied to a page only when the page is being rendered on a screen smaller than or equal to 800 pixels in width. A page of an RWD site typically counts multiple `LINK` elements: If no match is found, then no style sheet is applied to the page.

The CSS media query language is based on a pair of Boolean operators (and, not) and a few browser properties: `width`, `height`, `aspect-ratio`, `color`, and `orientation`. The most used of these properties is `width`, which refers to the width of the browser window. In CSS media queries, the `width` property can be decorated with a few prefixes: `device`, `min`, and `max`. The `device-width` property, in particular, refers to the physical width of the device.

Although RWD was not specifically designed for mobile scenarios, its inherent flexibility makes it more than suitable for mobile. From the perspective of an RWD-enabled site, a smartphone is simply a browser with a width of about 400 pixels. The site then detects the feature and adapts automatically.

Is your code secure?

Find out with HP Fortify on Demand.

Try Now



IN THIS ISSUE[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

This is the greatest strength of RWD. Unfortunately, it is important to note that it is also its major weakness.

Server-Side Responsive Sites

A site that can render well on a small screen without knowing it is running on some mobile device is a site that is not optimized for a mobile scenario. RWD and CSS media queries don't make it possible, for example, to switch the stylesheet to differentiate between an iPhone and an Android device. All that RWD can do is what is possible through CSS: hiding blocks or moving them around by leveraging the `float` CSS attribute.

When CSS is not enough to define the ideal mobile experience, then RWD is no longer an ideal choice. A website that just presents data — for example, a news portal — lends itself well to RWD. A more behavior-oriented site that implements workflows — for example, a booking site — requires ad hoc forms to provide an ideal experience. In this case, RWD alone is not sufficient.

The alternative to RWD is sniffing and analyzing the user agent (UA) string on the Web server. Once the user agent has been mapped to a device profile, the developer has all the information needed to intelligently serve ad hoc markup.

Does this mean that pages must be able to distinguish thousands of devices? No, the most common approach these days is to define a few classes of devices for which the site is then optimized. As an example, consider the following classes: smartphones, tablets, laptops, smart TVs. Then, use some professional tool to do the analysis of the user agent string and return known capabilities of each profiled device. Today, the de facto standard for UA-sniffing libraries is WURF.

A server-side approach to the design of mobile views doesn't mean you have to drop RWD or its growing ecosystem of related libraries such as Foundation and Bootstrap. Responsive Design Server Side (RESS) refers to adding a second dimension to the problem. You map the UA to a class of devices and define a view (and not simply a CSS) for each scenario. Defining a view, and not simply a CSS file as in RWD, gives you the opportunity to use the ideal layout for each case, size images appropriately, and minimize the markup being downloaded. The markup being downloaded, then, is just markup; as such, it can mimic RWD design principles and techniques.

Beyond Mobile

RWD represents a smart way to adapt sites to a variety of screen sizes. You should consider it for use in designing multi-device viewable websites. It's not the only solution, especially in cases where sites must be device specific, but it's an awful lot better than not accommodating mobile and portable devices at all.

— *Dino Esposito is a frequent contributor to Dr. Dobbs's and has written several books on designing and implementing Web applications.*

[Comment](#)

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

Supercharging ASP.NET Web Form Apps

With the introduction of the ASP.NET Web API in ASP.NET 4.5, it is now possible to add advanced client-side functionality to existing sites

By Charles Nurse

Most ASP.NET websites are built using the ASP.NET Web Forms model. It provides a simple, stateful abstraction over the stateless HTTP protocol, and focuses on server-side processing. Recently, websites have focused more on a client-side model with the use of JavaScript, AJAX, and REST-style services. With the introduction of the ASP.NET Web API in ASP.NET 4.5, it is now possible to add advanced client-side functionality to existing sites, while retaining many of the benefits of the Web Forms model.

A Simple Web Forms Application

To understand how to supercharge an existing Web Forms application with the ASP.NET Web API, let's first take a brief look at a Web Forms app that I'll use as an example. Figure 1 shows the list view of a simple Tasks application, which was created using the new Web Forms template shipped as part of Visual Studio 2012.

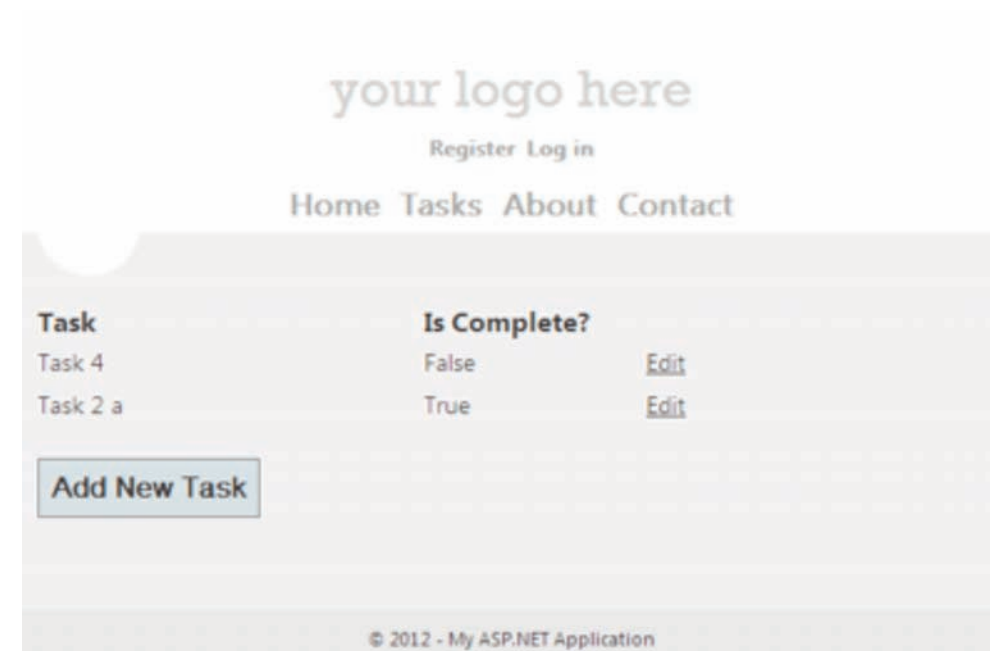


Figure 1: The Tasks page in list mode.

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

The Add New Task button and the Edit links in the Grid allow the user to add new tasks or edit the relevant existing tasks.

Assume that you have been given a new requirement: Instead of just displaying the “Is Complete” status of the task, the application should display a checkbox and allow users to update the status of the task (rather than having to click “Edit” and be redirected to the edit page).

You could fulfill this requirement by using a checkbox column in the ASP.NET `GridView` control and switching the `GridView` into Edit mode, but this would require the use of server post-backs and would trigger a complete refresh of the grid whenever a task is updated.

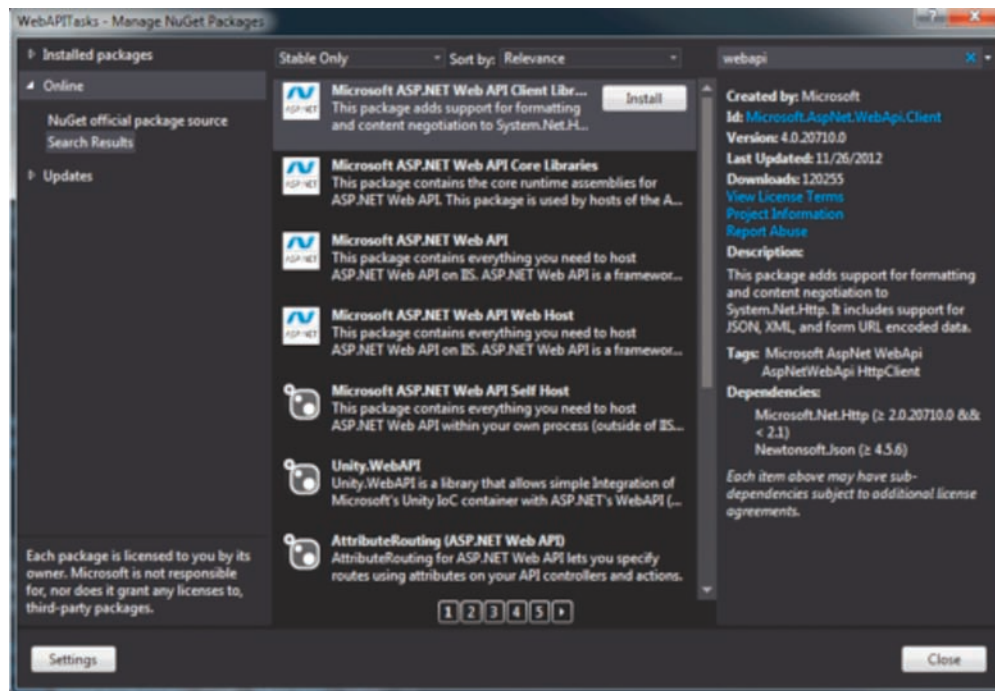


Figure 2: Adding the Web API using Nuget.

Instead, let’s use some JavaScript in the client and the ASP.NET Web API to demonstrate how to make the updates more modern and responsive.

Adding the ASP.NET Web API Using Nuget

ASP.NET Web API was released in conjunction with the .NET 4.5/Visual Studio 2012 release, but it doesn’t actually require .NET 4.5. Instead, it’s built against the earlier .NET 4.0 release. To add Web API support to an existing project, we can use Nuget to install the latest version, either via the website or via one of the integration points such as Manage Nuget Packages in Visual Studio (Figure 2).

This ensures you have the relevant components installed and will allow you to update those components if newer versions of the ASP.NET Web API are released. Note that, depending on what features of the ASP.NET Web API you want to use, you may still have to add some assembly references manually.

Adding a Web API Controller

The first step in adding our new feature using the Web API is to create a Web API Controller class (Figure 3). The Web API has some similarities to ASP.NET MVC in that the new class extends the `ApiController` base class and, by convention, is called `TasksController`. I will show

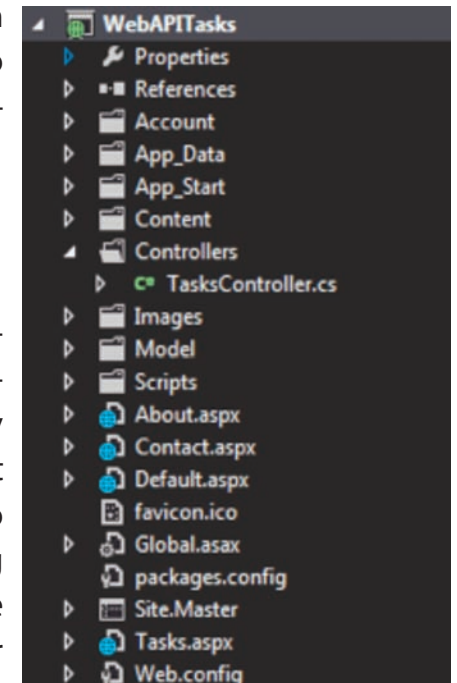


Figure 3: Adding a new Web API controller class.

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

how this convention works when we review the URI that needs to be called for our methods to be executed.

In the controller class, I will create a single method, `UpdateStatus`, that updates the task and changes the value of the `IsComplete` property, taking the `ID` as a single parameter (Listing One). As this is an update operation, I follow the convention of making the action method work with the `PUT` HTTP method (verb). I indicate that by using the `HttpPut` attribute on the method.

Listing One: The UpdateStatus Method.

```
[HttpPut]
public HttpResponseMessage UpdateStatus(int id)
{
    try
    {
        var db = new TasksContext();

        Task task = (from t in db.Tasks
                    where t.TaskID == id
                    select t).SingleOrDefault();

        if (task == null)
        {
            return Request.CreateResponse(
                HttpStatusCode.NotFound); ;
        }

        task.IsComplete = !task.IsComplete;

        db.SaveChanges();

        return Request.CreateResponse(HttpStatusCode.OK);
    }
    catch (Exception exc)
    {
        return Request.CreateErrorResponse(
            HttpStatusCode.InternalServerError, exc);
    }
}
```

In this simple example, I am using the Entity Framework for persistence, and the first few lines of this method retrieve the task to be updated. As long as I have a valid task (that is, the task is not null), I simply change the value of the `IsComplete` property and save the changes. This is fairly straightforward Entity Framework code.

If the update is successful, I create a new `HttpResponseMessage` with the `OK` status code (200), and this response is sent back to the caller. If there are exceptions, I create an `HttpResponseMessage` with the appropriate status code and, in the case of an Internal Server Error, details of the exception.

Adding Routing Information to Application_Start

Now that I have a Web API Controller class to route to, I need to set up the routes in the `Application_Start` event handler, so the Web API can execute the correct method. ASP.NET Web API routes are similar to ASP.NET MVC Routes. Two example routes are shown in Listing Two.

Listing Two: Adding ASP.NET Web API Routes to Application_Start.

```
RouteTable.Routes.MapHttpRoute(
    name: "DefaultApi",
    routeTemplate: "api/{controller}/{id}",
    defaults: new { id = RouteParameter.Optional }
);

RouteTable.Routes.MapHttpRoute(
    name: "TasksApi",
    routeTemplate: "api/{controller}/{action}/{id}",
    defaults: new { id = RouteParameter.Optional }
);
```

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

The first route is the default Web API route. Using this route, it's possible to use URIs like `api/Tasks/5`.

The `{controller}` place-holder routes control to the relevant Web API controller class. In the aforementioned example, it will be routed to the `TasksController`, and `api/products` would route to the `ProductsController`.

If used with the `GET` HTTP method (or verb), this URI is mapped to an action method that starts with the prefix `Get`, such as `GetTask` or `GetTaskById`, and the parameter (5) is passed as the ID parameter of the action method. Similarly, if used with the `DELETE` HTTP method, this URI is mapped to an action method that starts with the prefix `Delete`; such as `DeleteTask` or `DeleteTaskById`. Note that as long as the method has the correct prefix, control is routed to the correct action method regardless of the rest of the name. Issues arise only if the routing is ambiguous because there are two or more methods that start with the same HTTP method prefix: `GetTaskById` and `GetTaskByCategory`.

I could have followed this convention and used the `Put` prefix in the name of the action method I created in Listing One, but the action method does not contain any of these HTTP method prefixes. I can handle this by defining a second route to send control explicitly to a specific action. In this case, the `{action}` place-holder maps to the name of the action method, so the URI would need to be `api/Tasks/UpdateStatus/5`.

This URI will now route to the `UpdateStatus` method of the `TasksController`, passing the value 5 as the `taskId` parameter.

Adding JavaScript to Call the API

So, now I have an Web API method and I've defined the routing. The last step to hook up the new feature is to add JavaScript in the HTML to call the API.

To modify the `GridView` control to render an interactive checkbox instead of rendering `True` or `False` text, I can replace the `BoundField` in the `GridView` control markup with a `TemplateField`. This `TemplateField` has an `ItemTemplate` that contains an ASP.NET checkbox control as well as a hidden field (Listing Three).

Listing Three: The Updated GridView Markup

```
<asp:GridView runat="server" ID="tasksGrid"
  AutoGenerateColumns="False"
  DataKeyNames="TaskID"
  OnRowEditing="editTask">
  <Columns>
    <asp:BoundField DataField="TaskName"
      ItemStyle-Width="200px" HeaderText="Task"/>
    <asp:TemplateField ItemStyle-Width="100px"
      HeaderText="Is Complete?">
      <ItemTemplate>
        <asp:CheckBox ID="isCompleteCheck" runat="server"
          Checked='<%=Eval("IsComplete") %>' />
        <asp:HiddenField ID="taskIdField" runat="server"
          Value='<%=Eval("TaskId") %>' />
        </ItemTemplate>
      </asp:TemplateField>
    <asp:CommandField ShowEditButton="True" EditText="Edit" />
  </Columns>
</asp:GridView>
```

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

The `Checkbox` control is bound to the `IsComplete` property of the `Task` object, and the hidden field is bound to the `TaskId` property. The page now renders interactive checkboxes as shown in Figure 4.

I now have interactive checkboxes, and I can use jQuery to wire up the click event of the checkbox to call the new Web API method, as in Listing Four.

Listing Four: jQuery to Call the Web API UpdateStatus Method

```
(function ($) {
    $(document).ready(function () {
        $('input[name$="isCompleteCheck"]').click(function () {
            var $check = $(this);
            var $hidden =
                $check.parent().children('input[name$="taskIdField"]');
            var taskId = $hidden.val();

            $.ajax({
                type: "PUT",
                cache: false,
                url: 'api/Tasks/UpdateStatus/' +
                    taskId.toString()
            }).done(function (msg) {
                alert("Task Updated: " + msg);
            });
        });
    });
})(jQuery);
```

This is pretty standard jQuery. I wire up the click event of the checkboxes in the `GridView` by using the "Attribute ends with" selector `input[name$="isCompleteCheck"]`. In the same way, I then obtain

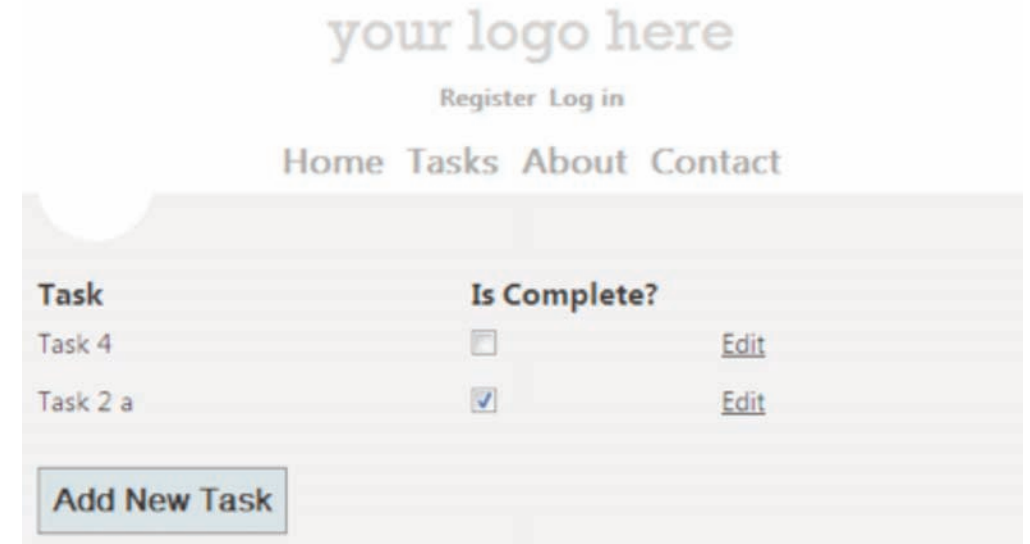


Figure 4: Updated Tasks page.

the relevant hidden field for that checkbox, and use the jQuery `val()` method to obtain the `taskId` of the selected checkbox. Finally, I use jQuery's `ajax` method to build the relevant call to the Web API method.

Note that I use the `PUT` HTTP method and create the URI discussed earlier in the article: `api/Tasks/UpdateStatus/{id}`. If the call is successful, the user is notified using a JavaScript alert dialog.

Taking It Further

I have only scratched the surface of what can be done with the Web API, but this example demonstrates that it is possible to add rich client-side functionality using the ASP.NET Web API with existing Web Forms applications. In the `GridView` control, I am still using server-side data binding to do the initial rendering of the list. But rather than do a full

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

post-back to the server just to update the `IsComplete` status, I added some simple JavaScript and a very simple AJAX call (combined with a Web API action method) to add client-side functionality. This provides a much more modern and responsive application.

This approach is an effective strategy for enhancing existing ASP.NET applications. I don't need to start from scratch and do a total rewrite in order to add some of these newer technologies. As emerging features are required, they can be added using jQuery and the ASP.NET Web API. Also, over time, legacy code could be adjusted to take advantage of this new approach. There are JavaScript libraries (such as KnockoutJS, Backbone, and Breeze) that can simplify the code as more and more of the logic is pushed to the client.

At DotNetNuke, we used this approach: We have a fundamental dependency on ASP.NET Web Forms and it is not feasible to rewrite everything to use ASP.NET MVC or the ASP.NET Web API. However, we were able to add a new Services Framework that is built on the Web API to provide this more modern functionality to third-party extension writers, as well as our own first-party extensions.

— *Charles Nurse has been a Senior Architect for DotNetNuke Corporation, the creators of the DotNetNuke Open Source Web Application Platform, since 2006. He is a Microsoft ASP.NET MVP and an ASPInsider.*

[Comment](#)


Instantly Search Terabytes Of Text

- 25+ fielded & full-text search options
- dtSearch's own file parsers **highlight hits** in popular file & email types
- Spider supports static & dynamic data
- APIs for .NET, Java, C++, SQL, etc.
- Win / Linux (64-bit & 32-bit)

"Lightning Fast" – *Redmond Mag*

"Covers all data sources" – *eWeek*

"Returns results in less than a second"
– *InfoWorld*

www.dtSearch.com

Fully-Functional Evaluations

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

TypeScript: The Lay of the Land

By creating a language that adds types to JavaScript, Microsoft made it easier to write complex Web apps, do compile-time syntax checking, and get better coding support in Visual Studio.

By **Gaston Hillar**

TypeScript is an open source superset of JavaScript under development by Microsoft. TypeScript adds support for optional static types, interfaces, classes, and modules to the JavaScript syntax, and translates this input into plain JavaScript. Thus, the resulting code runs in any browser or host. In this article, I introduce some of the most important TypeScript features and show how related tools use them to simplify program construction and maintenance, especially for complex applications. (An upcoming article will explore more subtle benefits of the language.)

JavaScript + Types

Most modern IDEs, including Visual Studio, have dramatically improved support for JavaScript. For example, Visual Studio 2012 provides good code-completion and navigation features for editing JavaScript code. When you work with an ASP.NET MVC 4 project, by default, the Scripts/_references.js file includes references to the files to which Visual

Studio will provide JavaScript IntelliSense support. This way, when you edit a JavaScript file, Visual Studio provides autocomplete for the most common libraries, just like jQuery.

However, when it is time to refactor, if you right-click on any variable, the context menu won't display the well-known Refactor | Rename... options (however, there are many third-party add-ins for Visual Studio 2012 that simplify this operation with JavaScript code). In addition, the autocomplete provides a "description of member variable" as the description for an element, which is not terribly helpful.

The lack of type information can produce runtime errors when you use unexpected types. Appropriate testing will detect the errors, but adding some type information would allow you to detect many errors when writing the code, and the editor would be able to provide you with accurate information about the error. That's the idea behind TypeScript.

If you have worked on an ASP.NET MVC application with a C# back end, you know that when you leave a C# file and start working on the

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

JavaScript pieces, Visual Studio loses a lot of features. For example, you don't have the dropdowns at the top of the editor that allow you to select a desired class and method to make it easier to navigate through the code. One of TypeScript's goals is to make the scripting experience in TypeScript similar in richness to the C# and C++ coding experience, without requiring a completely new scripting language.

As the number of JavaScript files and lines in an application grows, it becomes even more difficult to keep an easy-to-maintain structure. Thus, TypeScript adds support for classes with inheritance, modules, and interfaces to help in solving this problem. The syntax to define classes and modules is aligned with the ECMAScript 6 proposals.

Understanding TypeScript with the Playground

The easiest way to understand how TypeScript works is to dabble with the browser-based Playground, which uses the TypeScript language service in the code editor. The same TypeScript language service is available for Visual Studio 2012 with a plug-in. You can visit the TypeScript Playground (<http://www.typescriptlang.org/Playground/>) and start exploring TypeScript code without having to install any additional software. The Playground shows a split-screen editor with the TypeScript code on the left side and the generated JavaScript code on the right (Figure 1).

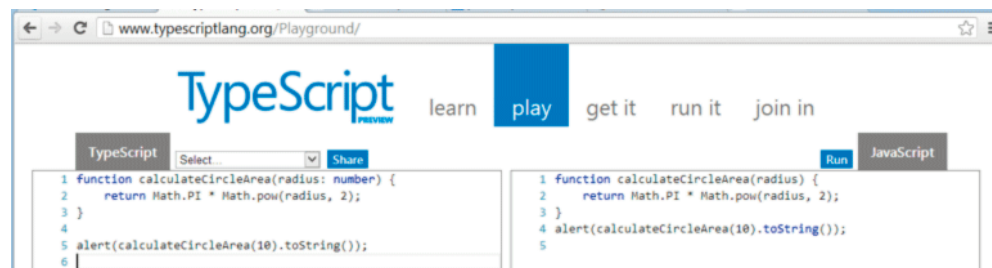


Figure 1: The TypeScript Playground displaying TypeScript code on the left and the generated JavaScript on the right.

TypeScript Playground

The following simple TypeScript code adds type information for the `radius` parameter (`number` type) of the `calculateCircleArea` function. The type inference mechanism will determine that the function returns a number, and therefore, it isn't necessary to specify the return type:

```
//// TypeScript code
function calculateCircleArea(radius: number) {
    return Math.PI * Math.pow(radius, 2);
}

alert(calculateCircleArea(10).toString());
```

As shown in Figure 1, the generated code is almost the identical: The JavaScript code just removes the type information added for the `radius` parameter, as shown in the following JavaScript lines:

```
//// JavaScript code generated by the TypeScript compiler
function calculateCircleArea(radius) {
    return Math.PI * Math.pow(radius, 2);
}

alert(calculateCircleArea(10).toString());
```

As you might guess, when TypeScript is translated to JavaScript, the type information is completely removed. The static types disappear at runtime, so there is zero cost for adding the types in TypeScript. However, because you provided type information for the `radius` parameter, the TypeScript editor is able to highlight errors if you supply the wrong parameter types for the `calculateCircleArea` function. For example, you can change the line that calls the `calculateCircleArea` function with the following line in the TypeScript editor:

```
alert(calculateCircleArea("10"));
```

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

The TypeScript editor knows that the `radius` parameters must be of a `number` type, but you supplied a `string` value. Therefore, the function name appears underlined in red. When you hover the mouse over the function name, the editor displays a tooltip with more details about the problem (see Figure 2): *“Supplied parameters do not match any signature of call target (radius: number) => number.”*

In this case, the code will execute without problems if you remove the type information. However, most of the time, you don't want to lose control of the conversions that happen under the hood, and you really want to receive a numeric value for your `calculateCircleArea` function. TypeScript provides static analysis based on both the type annotations you add and the structure of the code.

The same will happen if you change the line that calls the `calculateCircleArea` function with the following line in the TypeScript editor. In this case, the number of supplied parameters is wrong:

```
alert (calculateCircleArea(10, 20));
```

The TypeScript editor also provides autocompletion features. For example, if you enter `alert (calculate` and then press `Ctrl + Space`, the editor displays a dropdown menu with the `calculateCircleArea` function, its required parameter, and the returned type (see Figure 3).

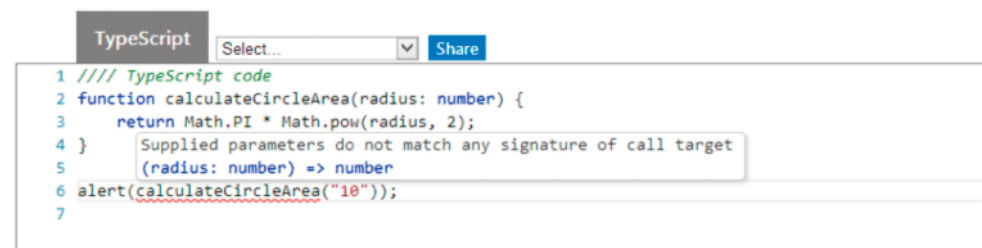


Figure 2: The TypeScript editor providing details about the error detected at design time.

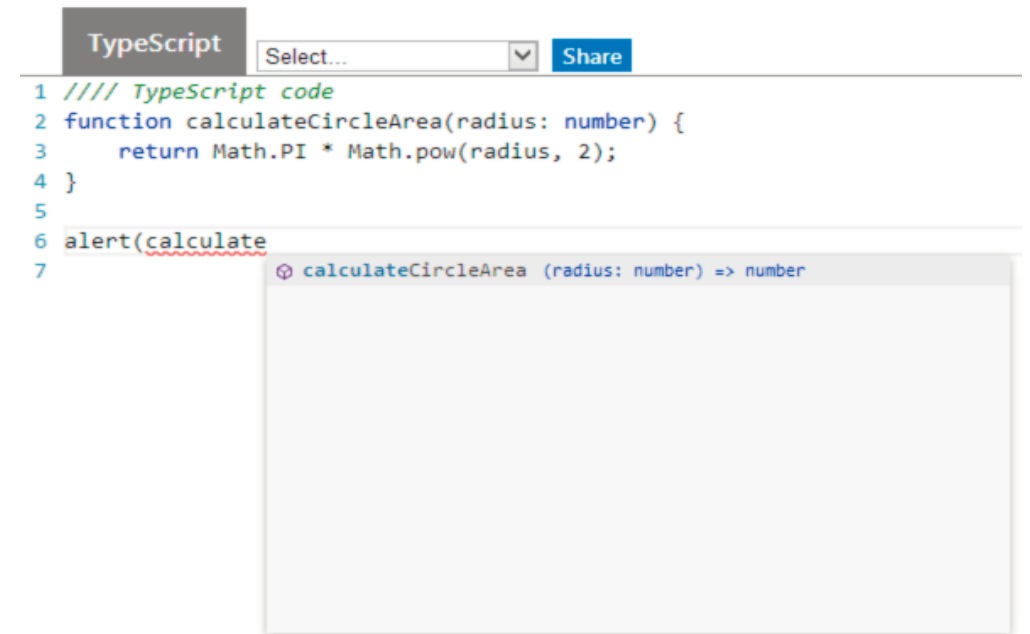


Figure 3: The TypeScript editor provides an autocomplete dropdown menu.

The editor will also highlight an error if you change the line that calls the `calculateCircleArea` function with the following line in the TypeScript editor (see Figure 4). In this case, the name `calculateCircleAreaWithRadius` doesn't exist in the current scope, and the code won't produce the expected results at runtime:

```
alert (calculateCircleAreaWithRadius(10));
```

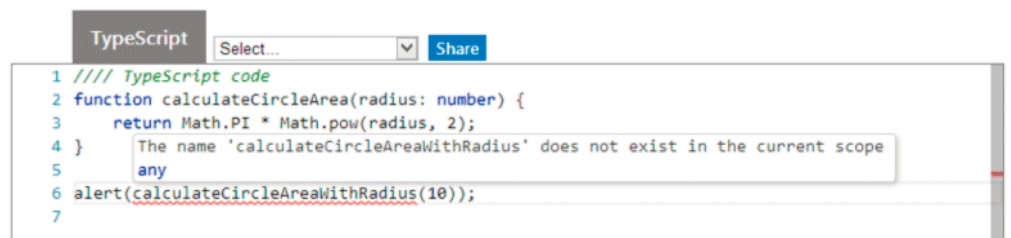


Figure 4: The TypeScript editor provides details about the error.

IN THIS ISSUE

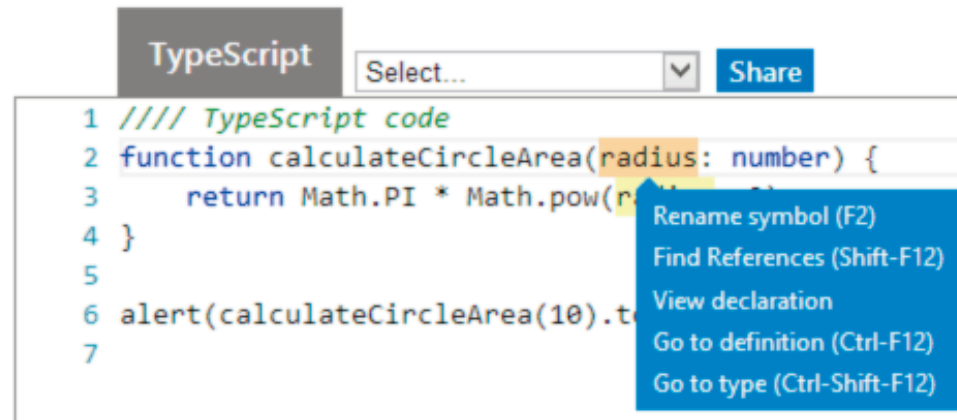
[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

Figure 5: The TypeScript editor displaying a context menu that allows you to rename a symbol.

If you click on the `radius` parameter and then move the mouse pointer, you will see a blue triangle appearing below `radius` and usage information about this symbol highlighted in the editor. If you click on the blue triangle, the editor displays a context menu with several options (Figure 5). Select “Rename symbol (F2)” and enter `rad` to change the name for `radius`. This way, you can easily rename a symbol or go to its definition. The same rename and definition navigation features are available when you install the TypeScript for Visual Studio 2012 plugin.

Keywords for Type Annotations

TypeScript supports the following keywords for type annotations to establish the intended contract of functions, variables, and properties:

- `any`: References the `Any` type. It represents any JavaScript value and it is a supertype of all types.
- `bool`: References the `Boolean` primitive type.
- `number`: References the `Number` primitive type that represents a double-precision 64-bit format IEEE 754 floating-point value.
- `null`: This literal references the only possible value of the `Null` type. The `Null` type is a subtype of all types (except `Void` and `Undefined`). It isn’t possible to explicitly reference the `Null` type and you can only use the `null` literal.
- `string`: References the `String` primitive type that represents sequences of characters stored as Unicode UTF-16 code units.
- `undefined`: This is the value given to uninitialized variables. The `Undefined` type is a subtype of all types. It isn’t possible to explicitly reference the `Undefined` type.
- `void`: This keyword references the `Void` type. It is the return type for functions that return no value and the only possible value of the `Void` type is `undefined`.

**VeriSign[®] SSL,
now from Symantec.**
More features. More protection.

Get more details now ▶



Confidence in a connected world.

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

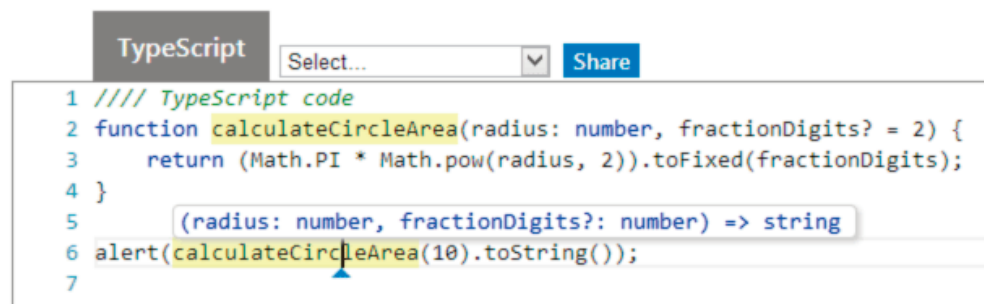
When you initialize variables or call functions to assign values to variables, you won't need to add explicit types because the type inference mechanism will generate the type information for both the TypeScript language service in the code editor and the compiler. Thus, you can detect errors related to types without having to use explicitly typed variables everywhere. For example, consider the following explicitly typed variable declaration:

```
var url: string = "www.drdobbs.com";
```

It isn't necessary to use the explicitly string typed variable declaration in order to enable type checking for TypeScript. TypeScript will know that the `url` variable is a `string` with just the following line:

```
var url = "www.drdobbs.com";
```

TypeScript allows you to easily specify optional parameters with default values. The following simple TypeScript code adds an optional `fractionDigits` parameter to the `calculateCircleArea` function by adding a question mark (?) to the parameter name. Notice that the default value for `fractionDigits` is 2. Thus, there is no need to explicitly specify the number type because TypeScript will infer it (Figure 6):



The screenshot shows the TypeScript editor interface. At the top, there is a 'TypeScript' label, a 'Select...' dropdown menu, and a 'Share' button. Below this, a code editor displays the following TypeScript code:

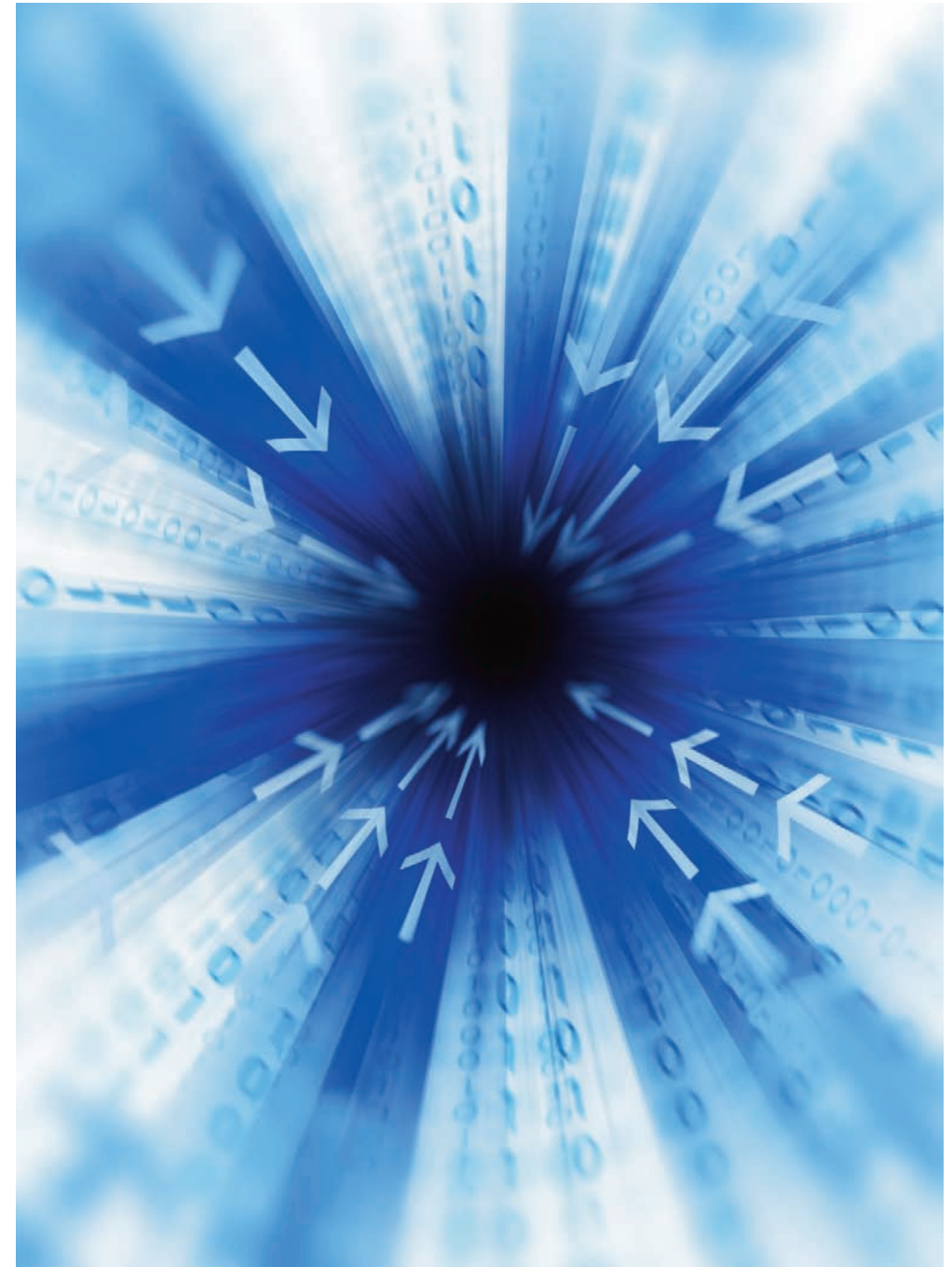
```

1 /// TypeScript code
2 function calculateCircleArea(radius: number, fractionDigits? = 2) {
3     return (Math.PI * Math.pow(radius, 2)).toFixed(fractionDigits);
4 }
5
6 alert(calculateCircleArea(10).toString());
7

```

A tooltip is visible over the `calculateCircleArea` function call on line 6. The tooltip contains the function signature: `(radius: number, fractionDigits?: number) => string`.

Figure 6: The TypeScript editor displaying a tooltip with the parameters and types for the `calculateCircleArea` function.



IN THIS ISSUE[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

```

////// TypeScript code
function calculateCircleArea(radius: number, fractionDigits?
= 2) {
    return (Math.PI * Math.pow(radius, 2)).toFixed
        (fractionDigits);
}

```

The following lines show the JavaScript code that the TypeScript compiler generates:

```

////// JavaScript code generated by the TypeScript compiler
function calculateCircleArea(radius, fractionDigits) {
    if (typeof fractionDigits === "undefined") {
        fractionDigits = 2; }
    return (Math.PI * Math.pow(radius, 2)).toFixed
        (fractionDigits);
}
alert (calculateCircleArea(10).toString());

```

In this case, the TypeScript compiler added the necessary JavaScript lines to check whether the type of `fractionDigits` is undefined. If the type is undefined, the code sets the default value of 2 to `fractionDigits`. This way, the call to the `calculateCircleArea` function with just the `radius` parameter specified works without problems and is a valid call. The syntax in TypeScript to define an optional parameter is easier to read than the JavaScript lines. In addition, TypeScript supports the object types, which will be further explored in a future article.

Interfaces

TypeScript supports interfaces as compile-time constructs that don't have a runtime representation; that is, they don't generate JavaScript code. You can use interfaces to declare a new named object-type and validate that required objects are passed as parameters and returned from functions.

[TYPESCRIPT]

Suppose you're developing a game that has aliens as the main characters. The following simple TypeScript code defines an `Alien` interface with the desired types for `name`, `angle`, `bonus`, and `rotationSpeed`. The `introduceAlien` function receives an `Alien` and uses both its name and bonus to display a message:

```

////// TypeScript code
interface Alien {
    name: string;
    angle: number;
    bonus: number;
    rotationSpeed: number;
}

function introduceAlien(alien: Alien) {
    return "I am " + alien.name + " and my bonus is " +
        alien.bonus.toString() + ".";
}

alert (introduceAlien({ name: "DrDobbsAlien", angle: 0,
    bonus: 5000, rotationSpeed: 5 }));

```

The following lines show the JavaScript code that the TypeScript compiler generates from this simple code. Notice that the interface doesn't generate JavaScript, so you just have the lines that define the `introduceAlien` function and the code that calls it:

```

////// JavaScript code generated by the TypeScript compiler
function introduceAlien(alien) {
    return "I am " + alien.name + " and my bonus is " +
        alien.bonus.toString() + ".";
}
alert (introduceAlien({
    name: "DrDobbsAlien",
    angle: 0,
    bonus: 5000,
    rotationSpeed: 5
})));

```

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

Now, if you were to enter `alert(introduceAlien({` and then press Ctrl + Space, the TypeScript editor would display a context menu with the four fields of the `Alien` interface (see Figure 7) because the editor knows that `introduceAlien` requires an `Alien` as defined in the interface. Thus, the interface enables you to reduce the number of runtime errors when functions require or return objects with certain fields.

TypeScript supports classes and inheritance, which allows you to create derived classes that specialize base classes. When you define classes, the TypeScript compiler generates JavaScript code for them,

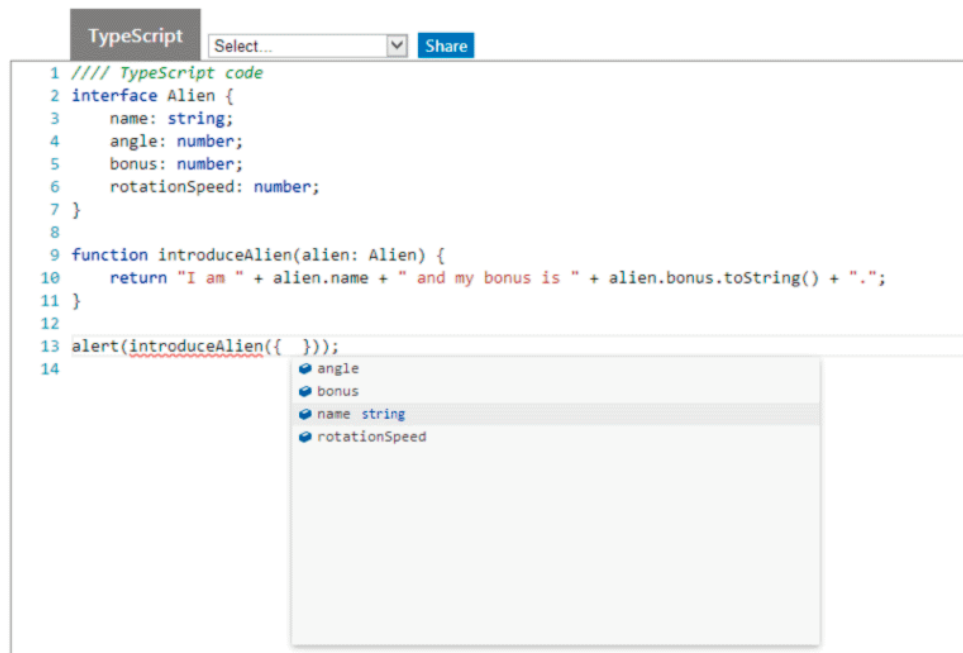


Figure 7: The TypeScript editor displaying a context menu with the four fields of the `Alien` interface.

but the member's accessibility disappears in JavaScript, as you will see in this example. TypeScript enforces private accessibility at compile-time, but there is no code that defines members as private at runtime. The following TypeScript code defines a `BadAlien` class that implements the previously introduced `Alien` interface. The code uses the previously shown `introduceAlien` function, which requires an `Alien` and, in this case, receives a `BadAlien` instance:

```

//// TypeScript code
class BadAlien implements Alien {
    damagePower: number;

    constructor(public name: string, public angle: number,
                public bonus: number, public rotationSpeed: number) {
        this.damagePower = 0.1 * bonus;
    }
}

interface Alien {
    name: string;
    angle: number;
    bonus: number;
    rotationSpeed: number;
}

function introduceAlien(alien: Alien) {
    return "I am " + alien.name + " and my bonus is " +
        alien.bonus.toString() + ".";
}

var firstBadAlien = new BadAlien("DrDobbsBadAlien", 0, 5000, 5);

alert(introduceAlien(firstBadAlien));

```

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

The following lines show the JavaScript code that the TypeScript compiler generates. Notice how the `BadAlien` class and its constructor are translated to functions:

```
//// JavaScript code generated by the TypeScript compiler
var BadAlien = (function () {
    function BadAlien(name, angle, bonus, rotationSpeed) {
        this.name = name;
        this.angle = angle;
        this.bonus = bonus;
        this.rotationSpeed = rotationSpeed;
        this.damagePower = 0.1 * bonus;
    }
    return BadAlien;
})();
function introduceAlien(alien) {
    return "I am " + alien.name + " and my bonus is " +
        alien.bonus.toString() + ".";
}
var firstBadAlien = new BadAlien("DrDobbsBadAlien", 0, 5000,
    5);
alert (introduceAlien (firstBadAlien));
```

[TYPESCRIPT]

Conclusion

In this article, I've explained the main benefits and features of TypeScript, and provided a brief overview of types, interfaces, and classes, focusing on the TypeScript Playground and the IntelliSense features that also work with Visual Studio 2012. In the next article in this series, I'll explain how you can work with TypeScript in Visual Studio 2012 with the latest versions of the plug-in and the compiler. In addition, I'll dive deeper into the support for classes, inheritance, properties, member accessibility, type assertions, arrow function expressions, modules, and automatically generated declarations that allow you to use all of TypeScript's handy features when consuming the most common JavaScript libraries.

— *Gaston Hillar is an expert in Windows-based development and is a frequent contributor to Dr. Dobb's.*

[Comment](#)

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

From the Vault

Building Scalable Web Architecture and Distributed Systems

Like most things in life, taking the time to plan ahead when building a Web service can help in the long run.

By **Kate Matsudaira**

Open source software has become a fundamental building block for some of the biggest websites. And as those websites have grown, best practices and guiding principles around their architectures have emerged. This article seeks to cover some of the key issues to consider when designing large websites, as well as some of the building blocks used to achieve these goals.

This article is largely focused on Web systems, although some of the material is applicable to other distributed systems as well.

Principles of Web Distributed Systems Design

What exactly does it mean to build and operate a scalable website or application? At a primitive level it's just connecting users with remote resources via the Internet — the part that makes it scalable is that the

resources, or access to those resources, are distributed across multiple servers.

Like most things in life, taking the time to plan ahead when building a Web service can help in the long run; understanding some of the considerations and tradeoffs behind big websites can result in smarter decisions at the creation of smaller websites. Below are some of the key principles that influence the design of large-scale Web systems:

- **Availability:** The uptime of a website is absolutely critical to the reputation and functionality of many companies. For some of the larger online retail sites, being unavailable for even minutes can result in thousands or millions of dollars in lost revenue, so designing their systems to be constantly available and resilient to

IN THIS ISSUE[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

failure is both a fundamental business and a technology requirement. High availability in distributed systems requires the careful consideration of redundancy for key components, rapid recovery in the event of partial system failures, and graceful degradation when problems occur.

- **Performance:** Website performance has become an important consideration for most sites. The speed of a website affects usage and user satisfaction, as well as search engine rankings, a factor that directly correlates to revenue and retention. As a result, creating a system that is optimized for fast responses and low latency is key.
- **Reliability:** A system needs to be reliable, such that a request for data will consistently return the same data. In the event the data changes or is updated, then that same request should return the new data. Users need to know that if something is written to the system, or stored, it will persist and can be relied on to be in place for future retrieval.
- **Scalability:** When it comes to any large distributed system, size is just one aspect of scale that needs to be considered. Just as important is the effort required to increase capacity to handle greater amounts of load, commonly referred to as the scalability of the system. Scalability can refer to many different parameters of the system: how much additional traffic can it handle, how easy is it to add more storage capacity, or even how many more transactions can be processed.
- **Manageability:** Designing a system that is easy to operate is another important consideration. The manageability of the system equates to the scalability of operations: maintenance and updates. Things to consider for manageability are the ease of diagnosing and understanding problems when they occur, ease of making updates or

modifications, and how simple the system is to operate (for example, does it routinely operate without failure or exceptions?).

- **Cost:** Cost is an important factor. This obviously can include hardware and software costs, but it is also important to consider other facets needed to deploy and maintain the system. The amount of developer time the system takes to build, the amount of operational effort required to run the system, and even the amount of training required should all be considered. Cost is the total cost of ownership.

Each of these principles provides the basis for decisions in designing a distributed Web architecture. However, they also can be at odds with one another, such that achieving one objective comes at the cost of another. A basic example: choosing to address capacity by simply adding more servers (scalability) can come at the price of manageability (you have to operate an additional server) and cost (the price of the servers).

When designing any sort of Web application it is important to consider these key principles, even if it is to acknowledge that a design may sacrifice one or more of them.

The Basics

When it comes to system architecture there are a few things to consider: What are the right pieces, how these pieces fit together, and what are the right tradeoffs. Investing in scaling before it is needed is generally not a smart business proposition; however, some forethought into the design can save substantial time and resources in the future.

This discussion is focused on some of the core factors that are central to almost all large Web applications: services, redundancy, partitions,

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

and handling failure. Each of these factors involves choices and compromises, particularly in the context of the principles described previously. In order to explain these in detail it is best to start with an example.

Example: Image Hosting Application

At some point you have probably posted an image online. For big sites that host and deliver lots of images, there are challenges in building an architecture that is cost-effective, highly available, and has low latency (fast retrieval).

Imagine a system where users are able to upload their images to a central server, and the images can be requested via a Web link or API, just like Flickr or Picasa. For the sake of simplicity, let's assume that this application has two key parts: the ability to upload (write) an image to the server, and the ability to query for an image. While we certainly want the upload to be efficient, we care most about having very fast delivery when someone requests an image (for example, images could be requested for a Web page or other application). This is very similar functionality to what a Web server or Content Delivery Network (CDN) edge server (a server CDN uses to store content in many locations so content is geographically/physically closer to users, resulting in faster performance) might provide.

Other important aspects of the system are:

- There is no limit to the number of images that will be stored, so storage scalability, in terms of image count, needs to be considered.
- There needs to be low latency for image downloads/requests.
- If a user uploads an image, the image should always be there (data reliability for images).

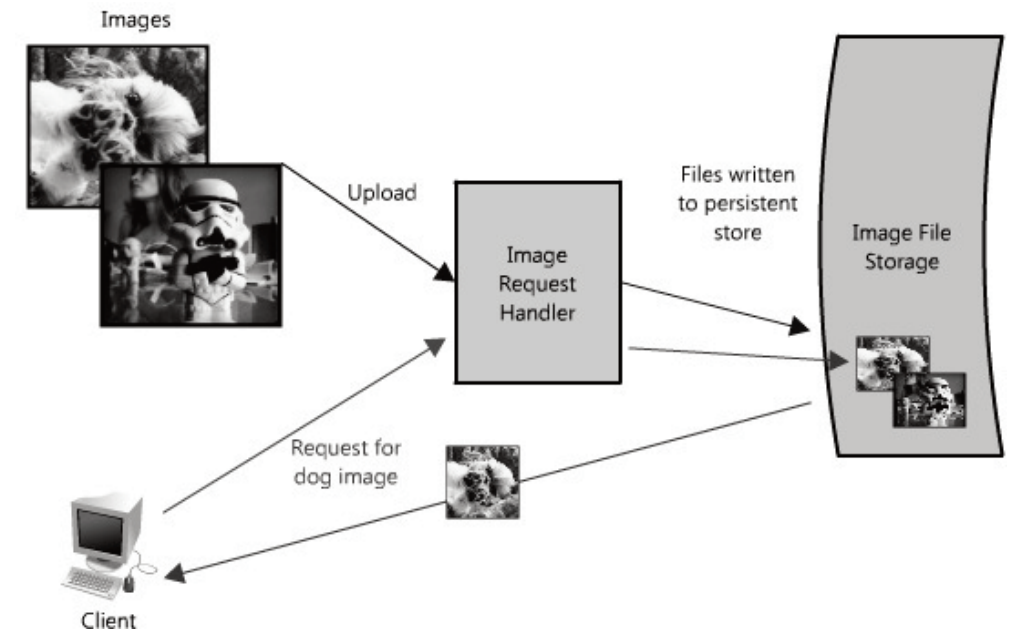


Figure 1: Simplified architecture diagram for image hosting application.

- The system should be easy to maintain (manageability).
- Since image hosting doesn't have high profit margins, the system needs to be cost-effective.

Figure 1 is a simplified diagram of the functionality. In this image hosting example, the system must be perceivably fast, its data stored reliably, and all of these attributes highly scalable. Building a small version of this application would be trivial and easily hosted on a single server; however, that would not be interesting for this discussion. Let's assume that we want to build something that could grow as big as Flickr.

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)**Services**

When considering scalable system design, it helps to decouple functionality and think about each part of the system as its own service with a clearly defined interface. In practice, systems designed in this way are said to have a Service-Oriented Architecture (SOA). For these types of systems, each service has its own distinct functional context, and interaction with anything outside of that context takes place through an abstract interface, typically the public-facing API of another service.

Deconstructing a system into a set of complementary services decouples the operation of those pieces from one another. This abstraction helps establish clear relationships between the service, its underlying environment, and the consumers of that service. Creating these clear delineations can help isolate problems, but also allows each piece to scale independently of one another. This sort of service-oriented design for systems is very similar to object-oriented design for programming.

In our example, all requests to upload and retrieve images are processed by the same server; however, as the system needs to scale, it makes sense to break out these two functions into their own services.

Fast-forward and assume that the service is in heavy use; such a scenario makes it easy to see how longer writes will impact the time it takes to read the images (since they two functions will be competing for shared resources). Depending on the architecture this effect can be substantial. Even if the upload and download speeds are the same (which is not true of most IP networks, since most are designed for at least a 3:1 download-speed:upload-speed ratio), read files will typically be read from cache, and writes will have to go to disk eventually (and perhaps be written several times in eventually consistent situations).

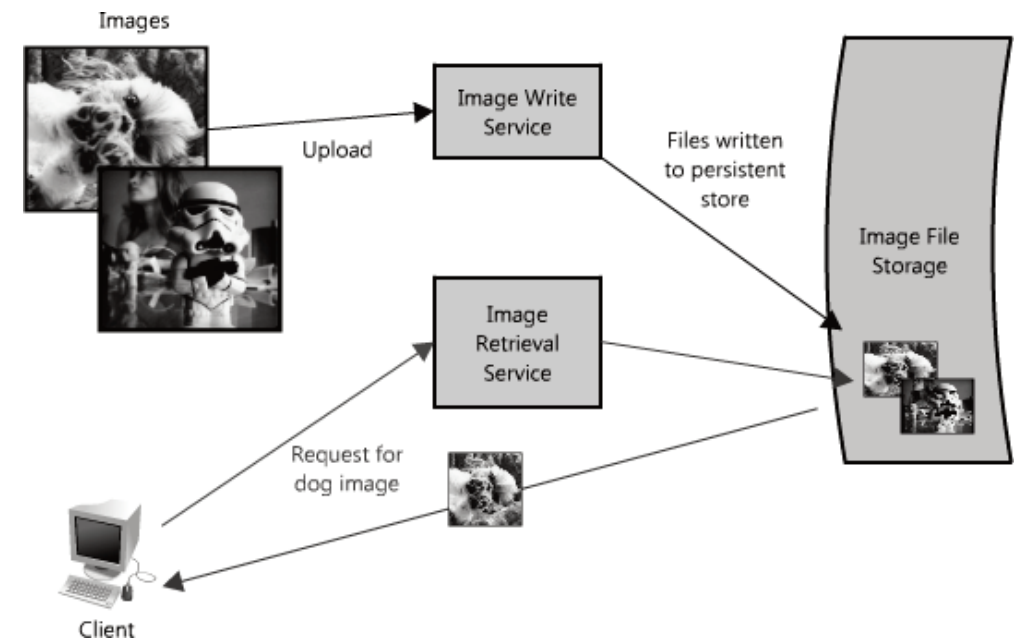


Figure 2: Splitting out reads and writes.

Even if everything is in memory or read from disks (like SSDs), database writes will almost always be slower than reads (Pole Position, an open source tool for DB benchmarking, <http://polepos.org/> and results <http://polepos.sourceforge.net/results/PolePositionClientServer.pdf>).

Another potential problem with this design is that a Web server like Apache or lighttpd typically has an upper limit on the number of simultaneous connections it can maintain (defaults are around 500, but can go much higher) and in high traffic, writes can quickly consume all of those. Since reads can be asynchronous, or take advantage of other performance optimizations like gzip compression or chunked transfer encoding, the Web server can switch serve reads faster and switch between clients quickly serving many more requests per sec-

IN THIS ISSUE[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

ond than the max number of connections (with Apache and max connections set to 500, it is not uncommon to serve several thousand read requests per second). Writes, on the other hand, tend to maintain an open connection for the duration for the upload, so uploading a 1MB file could take more than 1 second on most home networks, so that Web server could only handle 500 such simultaneous writes.

Planning for this sort of bottleneck makes a good case to split out reads and writes of images into their own services, shown in Figure 2. This allows us to scale each of them independently (since it is likely we will always do more reading than writing), but also helps clarify what is going on at each point. Finally, this separates future concerns, which would make it easier to troubleshoot and scale a problem like slow reads.

The advantage of this approach is that we are able to solve problems independently of one another — we don't have to worry about writing and retrieving new images in the same context. Both of these services still leverage the global corpus of images, but they are free to optimize their own performance with service-appropriate methods (for example, queuing up requests, or caching popular images — more on this below). And from a maintenance and cost perspective each service can scale independently as needed, which is great because if they were combined and intermingled, one could inadvertently impact the performance of the other as in the scenario discussed above.

Of course, the above example can work well when you have two different endpoints (in fact this is very similar to several cloud storage providers' implementations and Content Delivery Networks). There are lots of ways to address these types of bottlenecks though, and each has different tradeoffs.

For example, Flickr solves this read/write issue by distributing users across different shards such that each shard can only handle

a set number of users, and as users increase more shards are added to the cluster (see the presentation on Flickr's scaling, <http://is.gd/Jt394n>). In the first example it is easier to scale hardware based on actual usage (the number of reads and writes across the whole system), whereas Flickr scales with their user base (but forces the assumption of equal usage across users so there can be extra capacity). In the former an outage or issue with one of the services brings down functionality across the whole system (no-one can write files, for example), whereas an outage with one of Flickr's shards will only affect those users. In the first example it is easier to perform operations across the whole dataset — for example, updating the write service to include new metadata or searching across all image metadata — whereas with the Flickr architecture each shard would need to be updated or searched (or a search service would need to be created to collate that metadata — which is in fact what they do).

When it comes to these systems there is no right answer, but it helps to go back to the principles at the start of this article, determine the system needs (heavy reads or writes or both, level of concurrency, queries across the data set, ranges, sorts, etc.), benchmark different alternatives, understand how the system will fail, and have a solid plan for when failure happens.

Redundancy

In order to handle failure gracefully a Web architecture must have redundancy of its services and data. For example, if there is only one copy of a file stored on a single server, then losing that server means losing that file. Losing data is seldom a good thing, and a common way of handling it is to create multiple, or redundant, copies.

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

This same principle also applies to services. If there is a core piece of functionality for an application, ensuring that multiple copies or versions are running simultaneously can secure against the failure of a single node.

Creating redundancy in a system can remove single points of failure and provide a backup or spare functionality if needed in a crisis. For example, if there are two instances of the same service running in production, and one fails or degrades, the system can failover to the healthy copy. Failover can happen automatically or require manual intervention.

Another key part of service redundancy is creating a shared-nothing architecture. With this architecture, each node is able to operate independently of one another and there is no central “brain” managing state or coordinating activities for the other nodes. This helps a lot with scalability since new nodes can be added without special conditions or knowledge. However, and most importantly, there is no single point of failure in these systems, so they are much more resilient to failure.

For example, in our image server application, all images would have redundant copies on another piece of hardware somewhere (ideally in a different geographic location in the event of a catastrophe like an earthquake or fire in the data center), and the services to access the images would be redundant, all potentially servicing requests (see Figure 3). (Load balancers are a great way to make this possible, but there is more on that to come.)

Partitions

There may be very large data sets that are unable to fit on a single server. It may also be the case that an operation requires too many computing resources, diminishing performance and making it neces-

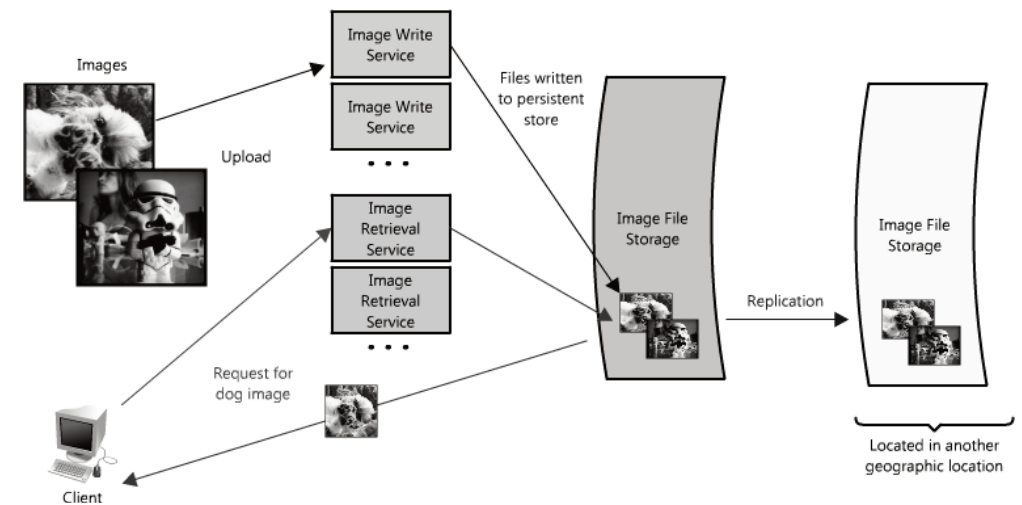


Figure 3: Image hosting application with redundancy.

sary to add capacity. In either case you have two choices: scale vertically or horizontally.

Scaling vertically means adding more resources to an individual server. So for a very large data set, this might mean adding more (or bigger) hard drives so a single server can contain the entire data set. In the case of the compute operation, this could mean moving the computation to a bigger server with a faster CPU or more memory. In each case, vertical scaling is accomplished by making the individual resource capable of handling more on its own.

To scale horizontally, on the other hand, is to add more nodes. In the case of the large data set, this might be a second server to store parts of the data set, and for the computing resource it would mean splitting the operation or load across some additional nodes. To take full advantage of horizontal scaling, it should be included as an intrinsic design principle of the system architecture, otherwise it can be quite cumbersome to modify and separate out the context to make this possible.

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

When it comes to horizontal scaling, one of the more common techniques is to break up your services into partitions, or shards. The partitions can be distributed such that each logical set of functionality is separate; this could be done by geographic boundaries, or by another criteria like non-paying versus paying users. The advantage of these schemes is that they provide a service or data store with added capacity.

In our image server example, it is possible that the single file server used to store images could be replaced by multiple file servers, each containing its own unique set of images. (See Figure 4.) Such an architecture would allow the system to fill each file server with images, adding additional servers as the disks become full. The design would require a naming scheme that tied an image's filename to the server containing it. An image's name could be formed from a consistent hashing scheme mapped across the servers. Or alternatively, each image could be assigned an incremental ID, so that when a client makes a request for an image, the image retrieval service only needs to maintain the range of IDs that are mapped to each of the servers (like an index).

Of course there are challenges distributing data or functionality across multiple servers. One of the key issues is data locality; in distributed systems the closer the data to the operation or point of computation, the better the performance of the system. Therefore it is potentially problematic to have data spread across multiple servers, as any time it is needed it may not be local, forcing the servers to perform a costly fetch of the required information across the network.

Another potential issue comes in the form of inconsistency. When there are different services reading and writing from a shared resource, potentially another service or data store, there is the chance for race conditions — where some data is supposed to be updated, but the read

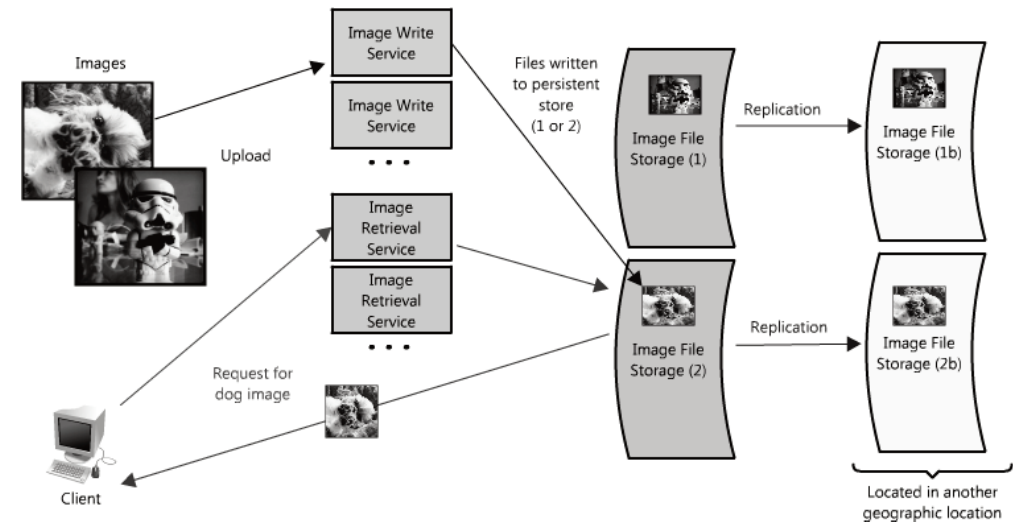


Figure 4: Image hosting application with redundancy and partitioning.

happens prior to the update — and in those cases the data is inconsistent. For example, in the image hosting scenario, a race condition could occur if one client sent a request to update the dog image with a new title, changing it from “Dog” to “Gizmo,” but at the same time another client was reading the image. In that circumstance it is unclear which title, “Dog” or “Gizmo,” would be the one received by the second client.

There are certainly some obstacles associated with partitioning data, but partitioning allows each problem to be split — by data, load, usage patterns, etc. — into manageable chunks. This can help with scalability and manageability, but is not without risk. There are lots of ways to mitigate risk and handle failures; however, in the interest of brevity they are not covered in this article. If you are interested in reading more, you can check out my blog post on fault tolerance and monitoring at <http://is.gd/u1BMPE>.

IN THIS ISSUE

- [Editorial >>](#)
- [ASP.NET >>](#)
- [TypeScript >>](#)
- [Scalable Web >>](#)
- [Links >>](#)
- [Table of Contents >>](#)



Figure 5: Simple web applications.

The Building Blocks of Fast and Scalable Data Access

Having covered some of the core considerations in designing distributed systems, let's now talk about the hard part: scaling access to the data.

Most simple Web applications, for example, LAMP stack applications, look something like Figure 5. As they grow, there are two main challenges: scaling access to the app server and to the database. In a highly scalable application design, the app (or Web) server is typically minimized and often embodies a shared-nothing architecture. This makes the app server layer of the system horizontally scalable. As a result of this design, the heavy lifting is pushed down the stack to the database server and supporting services; it's at this layer where the real scaling and performance challenges come into play.

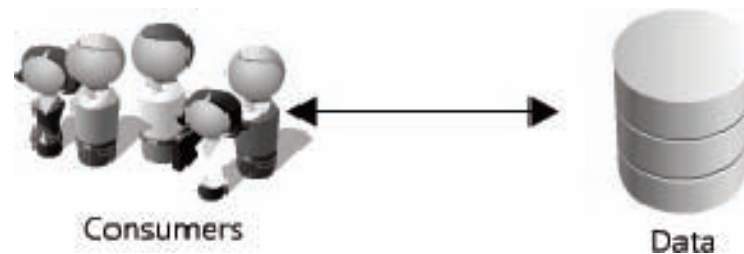


Figure 6: Oversimplified web application.

The rest of this article is devoted to some of the more common strategies and methods for making these types of services fast and scalable by providing fast access to data.

Most systems can be oversimplified to Figure 6. This is a great place to start. If you have a lot of data, you want fast and easy access, like keeping a stash of candy in the top drawer of your desk. Though overly simplified, the previous statement hints at two hard problems: scalability of storage and fast access of data.

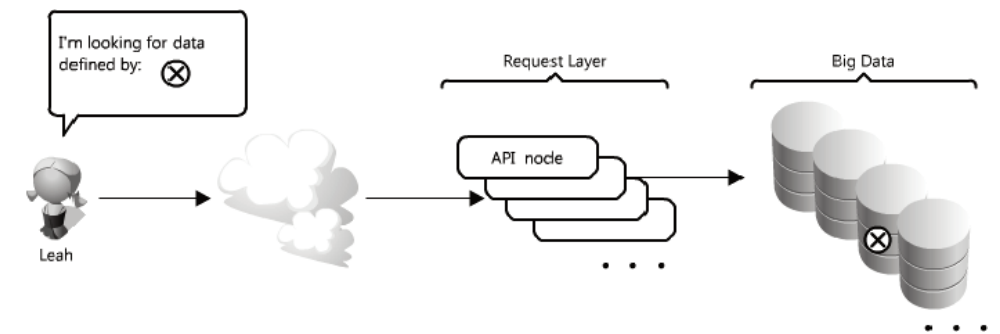


Figure 7: Accessing specific data.

For the sake of this discussion, let's assume you have many terabytes (TB) of data and you want to allow users to access small portions of that data at random. (See Figure 7.) This is similar to locating an image file somewhere on the file server in the image application example.

This is particularly challenging because it can be very costly to load TBs of data into memory; this directly translates to disk IO. Reading from disk is many times slower than from memory — memory access is as fast as Chuck Norris, whereas disk access is slower than the line at the DMV. This speed difference really adds up for large data sets; in real numbers memory access is as little as 6 times faster for sequential reads, or 100,000 times faster for random reads, than reading from disk

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

(see “The Pathologies of Big Data,” <http://is.gd/MI5zux>). Moreover, even with unique IDs, solving the problem of knowing where to find that little bit of data can be an arduous task. It’s like trying to get that last Jolly Rancher from your candy stash without looking.

Thankfully there are many options that you can employ to make this easier; four of the more important ones are caches, proxies, indexes and load balancers. The rest of this article discusses how each of these concepts can be used to make data access a lot faster.

Caches

Caches take advantage of the locality of reference principle: recently requested data is likely to be requested again. They are used in almost every layer of computing: hardware, operating systems, Web browsers, Web applications and more. A cache is like short-term memory: it has a limited amount of space, but is typically faster than the original data source and contains the most recently accessed items. Caches can exist at all levels in architecture, but are often found at the level nearest to the front end, where they are implemented to return data quickly without taxing downstream levels.

How can a cache be used to make your data access faster in our API example? In this case, there are a couple of places you can insert a cache. One option is to insert a cache on your request layer node, as in Figure 8.

Placing a cache directly on a request layer node enables the local storage of response data. Each time a request is made to the service, the node will quickly return local, cached data if it exists. If it is not in the cache, the request node will query the data from disk. The cache on one request layer node could also be located both in memory (which is very fast) and on the node’s local disk (faster than going to

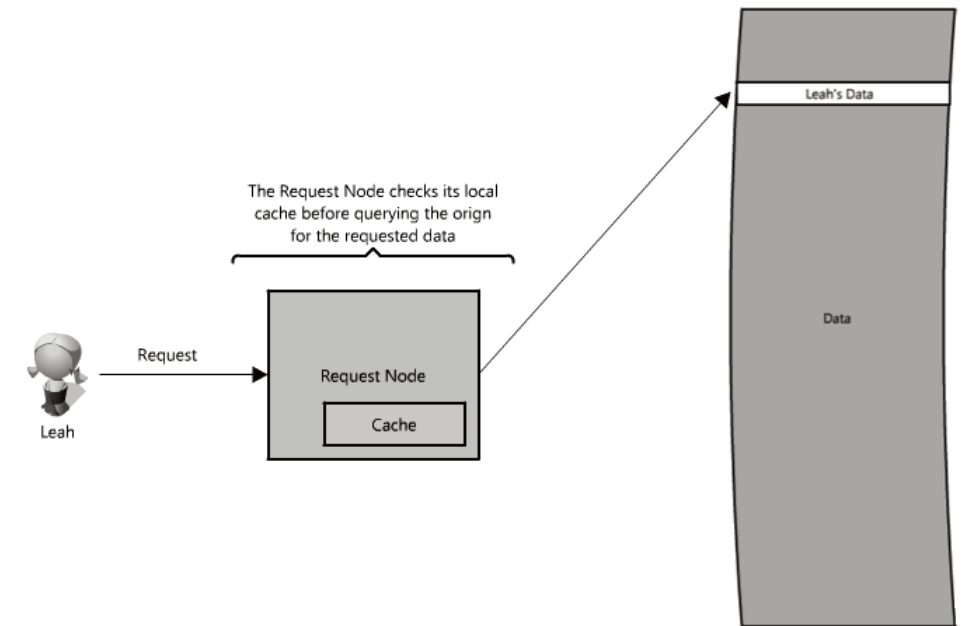


Figure 8: Inserting a cache on your request layer node.

network storage).

What happens when you expand this to many nodes? As you can see in Figure 9, if the request layer is expanded to multiple nodes, it’s still quite possible to have each node host its own cache. However, if your load balancer randomly distributes requests across the nodes, the same request will go to different nodes, thus increasing cache misses. Two choices for overcoming this hurdle are global caches and distributed caches.

Global Cache

A global cache is just as it sounds: all the nodes use the same single cache space. This involves adding a server, or file store of some sort, faster than your original store and accessible by all the request layer

IN THIS ISSUE

- [Editorial >>](#)
- [ASP.NET >>](#)
- [TypeScript >>](#)
- [Scalable Web >>](#)
- [Links >>](#)
- [Table of Contents >>](#)

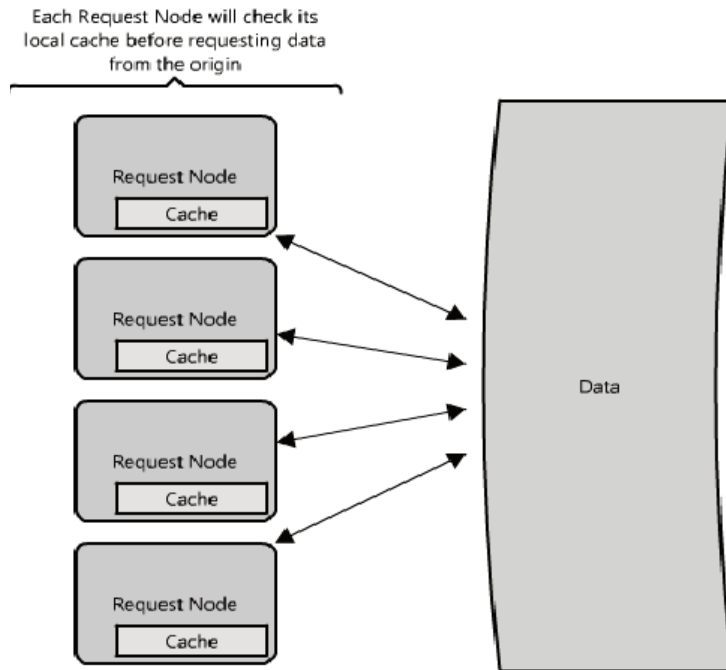


Figure 9: Multiple caches..

nodes. Each of the request nodes queries the cache in the same way it would a local one. This kind of caching scheme can get a bit complicated because it is very easy to overwhelm a single cache as the number of clients and requests increase, but is very effective in some architectures (particularly ones with specialized hardware that make this global cache very fast, or that have a fixed dataset that needs to be cached).

There are two common forms of global caches depicted in the diagrams. In Figure 10, when a cached response is not found in the cache, the cache itself becomes responsible for retrieving the missing piece of data from the underlying store. In Figure 11 it is the responsibility of request nodes to retrieve any data that is not found in the cache.

The majority of applications leveraging global caches tend to use the first type, where the cache itself manages eviction and fetching data

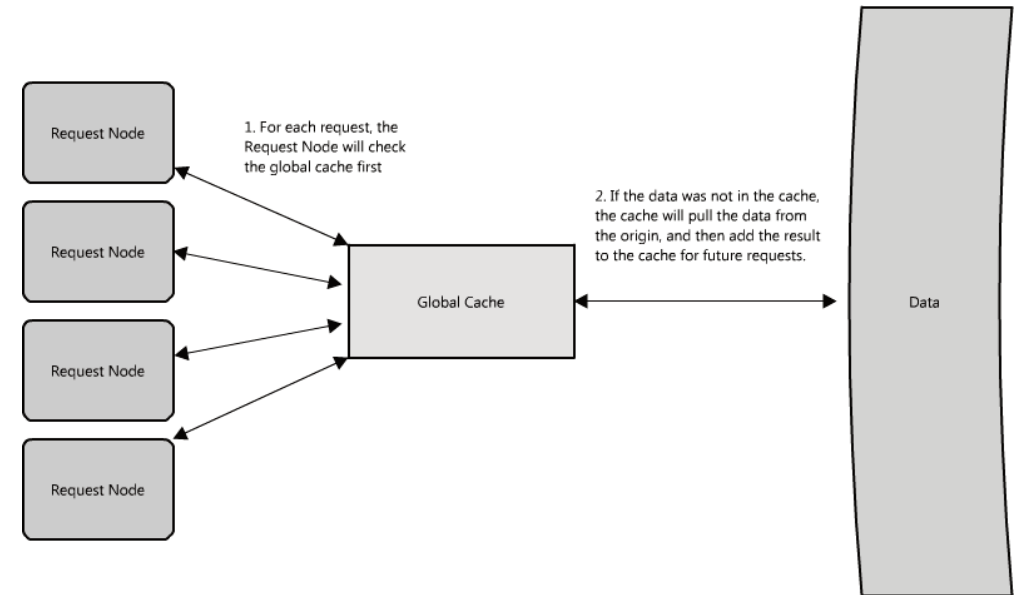


Figure 10: Global cache where cache is responsible for retrieval.

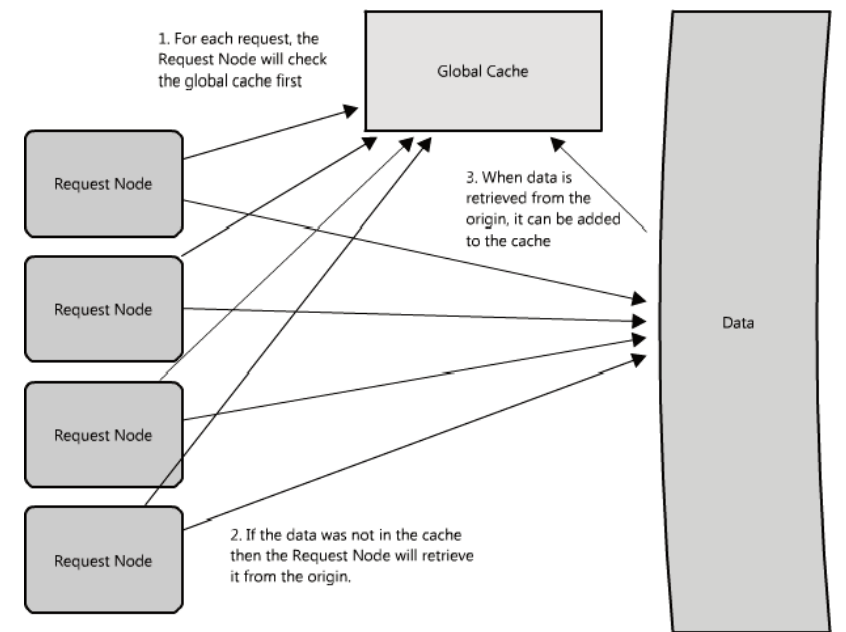


Figure 11: Global cache where request nodes are responsible for retrieval.

IN THIS ISSUE

[Editorial >>](#)

[ASP.NET >>](#)

[TypeScript >>](#)

[Scalable Web >>](#)

[Links >>](#)

[Table of Contents >>](#)

to prevent a flood of requests for the same data from the clients. However, there are some cases where the second implementation makes more sense. For example, if the cache is being used for very large files, a low cache hit percentage would cause the cache buffer to become overwhelmed with cache misses; in this situation it helps to have a large percentage of the total data set (or hot data set) in the cache. Another example is an architecture where the files stored in the cache are static and shouldn't be evicted. (This could be because of application requirements around that data latency — certain pieces of data might need to be very fast for large data sets — where the application logic understands the eviction strategy or hot spots better than the cache.)

Distributed Cache

In a distributed cache (Figure 12), each of its nodes own part of the cached data, so if a refrigerator acts as a cache to the grocery store, a distributed cache is like putting your food in several locations — your fridge, cupboards, and lunch box — convenient locations for retrieving snacks from, without a trip to the store. Typically the cache is divided up using a consistent hashing function, such that if a request node is looking for a certain piece of data it can quickly know where to look within the distributed cache to determine if that data is available. In this case, each node has a small piece of the cache, and will then send a request to another node for the data before going to the origin. Therefore, one of the advantages of a distributed cache is the increased cache space that can be had just by adding nodes to the request pool.

A disadvantage of distributed caching is remedying a missing node. Some distributed caches get around this by storing multiple copies of the data on different nodes; however, you can imagine how this logic

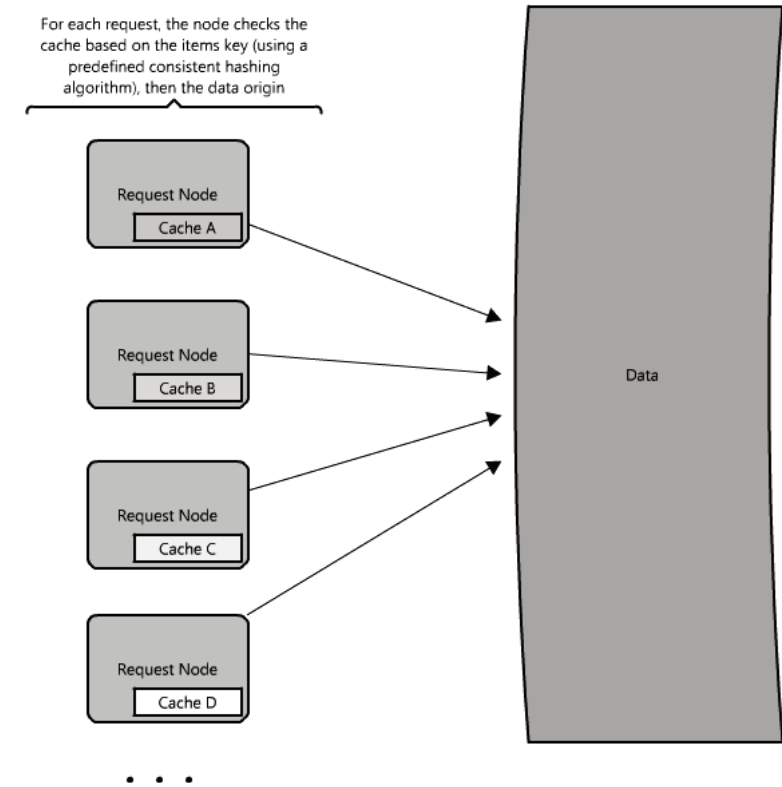


Figure 12: Distributed cache.

can get complicated quickly, especially when you add or remove nodes from the request layer. Although even if a node disappears and part of the cache is lost, the requests will just pull from the origin — so it isn't necessarily catastrophic!

The great thing about caches is that they usually make things much faster (implemented correctly, of course!) The methodology you choose just allows you to make it faster for even more requests. However, all this caching comes at the cost of having to maintain additional storage space, typically in the form of expensive memory; nothing is

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

free. Caches are wonderful for making things generally faster, and moreover provide system functionality under high load conditions when otherwise there would be complete service degradation.

One example of a popular open source cache is Memcached (<http://memcached.org/>) (which can work both as a local cache and distributed cache); however, there are many other options (including many language- or framework-specific options).

Memcached is used in many large websites, and even though it can be very powerful, it is simply an in-memory key value store, optimized for arbitrary data storage and fast lookups ($O(1)$).

Facebook uses several different types of caching to obtain their site performance (see “Facebook caching and performance” at <http://sizzo.org/wp/talks>). They use `$GLOBALS` and APC caching at the language level (provided in PHP at the cost of a function call) which helps make intermediate function calls and results much faster. (Most languages have these types of libraries to improve Web page performance and they should almost always be used.) Facebook then use a global cache that is distributed across many servers (see “Scaling memcached at Facebook” at <http://is.gd/73yFDm>), such that one function call accessing the cache could make many requests in parallel for data stored on different Memcached servers. This allows them to get much higher performance and throughput for their user profile data, and have one central place to update data (which is important, since cache invalidation and maintaining consistency can be challenging when you are running thousands of servers).

Now let’s talk about what to do when the data isn’t in the cache...

Proxies

At a basic level, a proxy server is an intermediate piece of hardware/software that receives requests from clients and relays them to the



Figure 13: Proxy server.

backend origin servers. Typically, proxies are used to filter requests, log requests, or sometimes transform requests (by adding/removing headers, encrypting/decrypting, or compression).

Proxies are also immensely helpful when coordinating requests from multiple servers, providing opportunities to optimize request traffic from a system-wide perspective. One way to use a proxy to speed up data access is to collapse the same (or similar) requests together into one request, and then return the single result to the requesting clients. This is known as collapsed forwarding.

Imagine there is a request for the same data (let’s call it `littleB`) across several nodes, and that piece of data is not in the cache. If that request is routed through the proxy, then all of those requests can be collapsed into one, which means we only have to read `littleB` off disk once. (See Figure 14.) There is some cost associated with this design, since each request can have slightly higher latency, and some requests may be slightly delayed to be grouped with similar ones. But it will improve performance in high load situations, particularly when that same data is requested over and over. This is similar to a cache, but instead of storing the data/document like a cache, it is optimizing the requests or calls for those documents and acting as a proxy for those clients.

IN THIS ISSUE

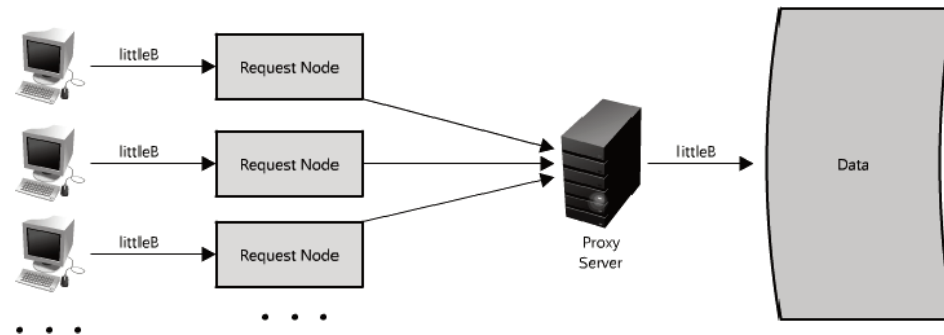
[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

Figure 14: Using a proxy server to collapse requests.

In a LAN proxy, for example, the clients do not need their own IPs to connect to the Internet, and the LAN will collapse calls from the clients for the same content. It is easy to get confused here though, since many proxies are also caches (as it is a very logical place to put a cache), but not all caches act as proxies.

Another great way to use the proxy is to not just collapse requests for the same data, but also to collapse requests for data that is spatially close together in the origin store (consecutively on disk). Employing such a strategy maximizes data locality for the requests, which can result in decreased request latency.

For example, let's say a bunch of nodes request parts of B: `partB1`, `partB2`, etc. We can set up our proxy to recognize the spatial locality of the individual requests, collapsing them into a single request and returning only `bigB`, greatly minimizing the reads from the data origin. (See Figure 15.) This can make a really big difference in request time when you are randomly accessing across TBs of data! Proxies are especially helpful under high load situations, or when you have limited caching, since they can essentially batch several requests into one.

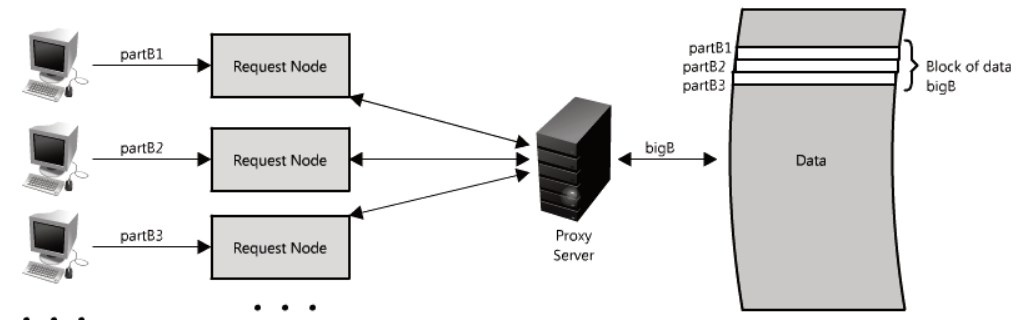


Figure 15: Using a proxy to collapse requests for data that is spatially close together.

It is worth noting that you can use proxies and caches together, but generally it is best to put the cache in front of the proxy, for the same reason that it is best to let the faster runners start first in a crowded marathon race. This is because the cache is serving data from memory, it is very fast, and it doesn't mind multiple requests for the same result. But if the cache was located on the other side of the proxy server, then there would be additional latency with every request before the cache, and this could hinder performance.

If you are looking at adding a proxy to your systems, there are many options to consider; Squid (<http://www.squid-cache.org/>) and Varnish (<https://www.varnish-cache.org/>) have both been road tested and are widely used in many production websites.

These proxy solutions offer many optimizations to make the most of client-server communication. Installing one of these as a reverse proxy (explained in the load balancer section) at the Web server layer can improve Web server performance considerably, reducing the amount of work required to handle incoming client requests.

IN THIS ISSUE

- [Editorial >>](#)
- [ASP.NET >>](#)
- [TypeScript >>](#)
- [Scalable Web >>](#)
- [Links >>](#)
- [Table of Contents >>](#)

Indexes

Using an index to access your data quickly is a well-known strategy for optimizing data access performance; probably the most well known when it comes to databases. An index makes the trade-offs of increased storage overhead and slower writes (since you must both write the data and update the index) for the benefit of faster reads.

Just as to a traditional relational data store, you can also apply this concept to larger data sets. The trick with indexes is you must carefully consider how users will access your data. In the case of data sets that are many TBs in size, but with very small payloads (e.g., 1 KB), indexes are a necessity for optimizing data access. Finding a small payload in such a large data set can be a real challenge since you can't possibly iterate over that much data in any reasonable time. Furthermore, it is very likely that such a large data set is spread over several (or many!) physical devices — this means you need some way to find the correct physical location of the desired data. Indexes are the best way to do this.

An index can be used like a table of contents that directs you to the location where your data lives. For example, let's say you are looking for a piece of data, part 2 of B — how will you know where to find it? If you have an index that is sorted by data type — say data A, B, C — it would tell you the location of data B at the origin. Then you just have to seek to that location and read the part of B you want. (See Figure 16.)

These indexes are often stored in memory, or somewhere very local to the incoming client request. Berkeley DBs (BDBs) and tree-like data structures are commonly used to store data in ordered lists, ideal for access with an index,

Often there are many layers of indexes that serve as a map, moving you from one location to the next, and so forth, until you get the specific piece of data you want. (See Figure 17.)

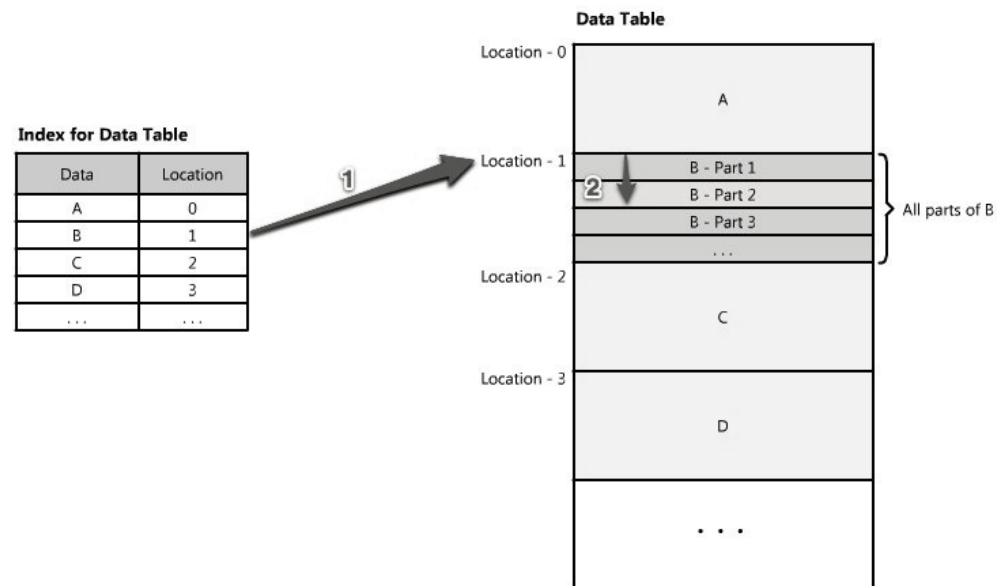


Figure 16: Indexes.

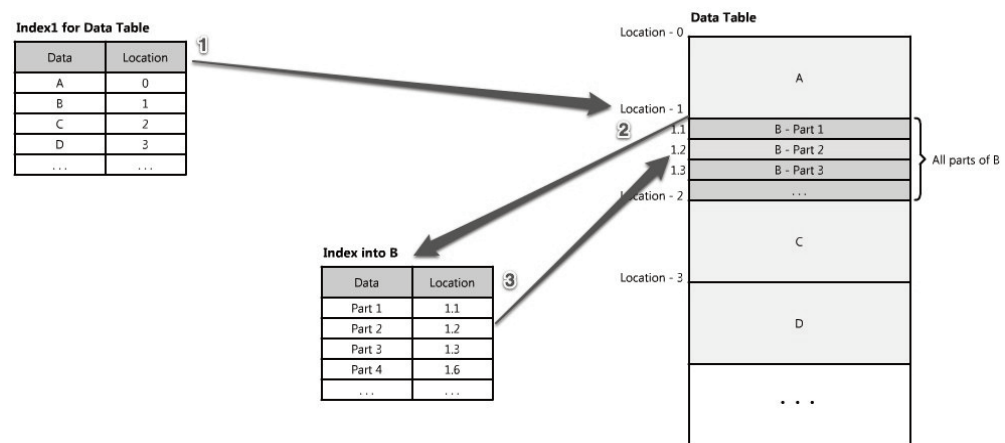


Figure 17: Many layers of indexes.

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

Indexes can also be used to create several different views of the same data. For large data sets, this is a great way to define different filters and sorts without resorting to creating many additional copies of the data.

For example, imagine that the image hosting system from earlier is actually hosting images of book pages, and the service allows client queries across the text in those images, searching all the book content about a topic, in the same way search engines allow you to search HTML content. In this case, all those book images take many, many servers to store the files, and finding one page to render to the user can be a bit involved. First, inverse indexes to query for arbitrary words and word tuples need to be easily accessible; then there is the challenge of navigating to the exact page and location within that book, and retrieving the right image for the results. So in this case the inverted index would map to a location (such as book B), and then B may contain an index with all the words, locations and number of occurrences in each part.

An inverted index, which could represent `Index1` in the Figure 17, might look something like the following — each word or tuple of words provide an index of what books contain them.

Word (s)	Book (s)
being awesome	Book B, Book C, Book D
always	Book C, Book F
believe	Book B

The intermediate index would look similar but would contain just the words, location, and information for book B. This nested index architecture allows each of these indexes to take up less space than if all of that info had to be stored into one big inverted index.

And this is key in large-scale systems because even compressed, these indexes can get quite big and expensive to store. In this system if we assume we have a lot of the books in the world — 100,000,000 (see Inside Google Books blog post at <http://is.gd/eG0ZcE>) — and that each book is only 10 pages long (to make the math easier), with 250 words per page, that means there are 250 billion words. If we assume an average of 5 characters per word, and each character takes 8 bits (or 1 byte, even though some characters are 2 bytes), so 5 bytes per word, then an index containing only each word once is over a terabyte of storage. So you can see creating indexes that have a lot of other information like tuples of words, locations for the data, and counts of occurrences, can add up very quickly.

Creating these intermediate indexes and representing the data in smaller sections makes big data problems tractable. Data can be spread across many servers and still accessed quickly. Indexes are a cornerstone of information retrieval, and the basis for today's modern search engines. Of course, this discussion only scratched the surface, and there is a lot of research being done on how to make indexes smaller, faster, contain more information (like relevancy), and update seamlessly. (There are some manageability challenges with race conditions, and with the sheer number of updates required to add new data or change existing data, particularly in the event where relevancy or scoring is involved).

Being able to find your data quickly and easily is important; indexes are an effective and simple tool to achieve this.

Load Balancers

Finally, another critical piece of any distributed system is a load balancer. Load balancers are a principal part of any architecture, as their

IN THIS ISSUE

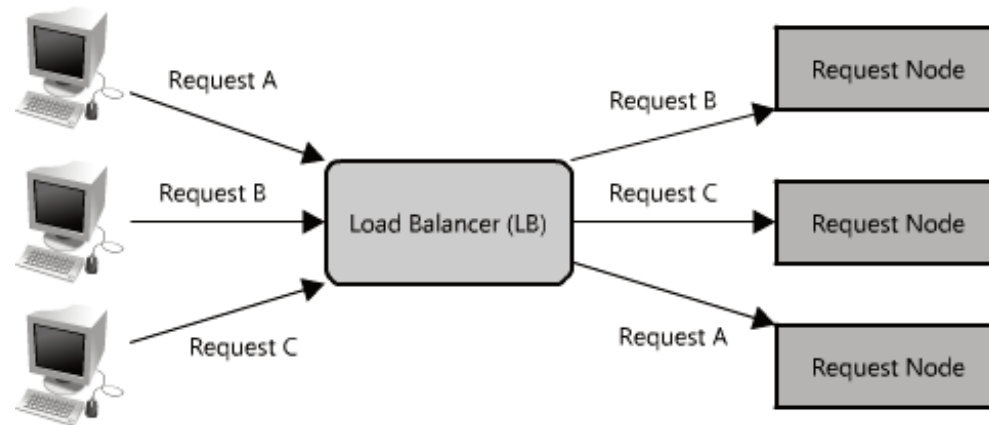
[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

Figure 18: Load balancer.

role is to distribute load across a set of nodes responsible for servicing requests. This allows multiple nodes to transparently service the same function in a system. (See Figure 18.) Their main purpose is to handle a lot of simultaneous connections and route those connections to one of the request nodes, allowing the system to scale to service more requests by just adding nodes.

There are many different algorithms that can be used to service requests, including picking a random node, round robin, or even selecting the node based on certain criteria, such as memory or CPU utilization. Load balancers can be implemented as software or hardware appliances. One open source software load balancer that has received wide adoption is HAProxy; see <http://haproxy.1wt.eu/>.

In a distributed system, load balancers are often found at the very front of the system, such that all incoming requests are routed accordingly. In a complex distributed system, it is not uncommon for a request to be routed to multiple load balancers as shown in Figure 19.

Like proxies, some load balancers can also route a request differently depending on the type of request it is. (Technically these are also known as reverse proxies.)

One of the challenges with load balancers is managing user-session-specific data. In an e-commerce site, when you only have one client it is very easy to allow users to put things in their shopping cart and persist those contents between visits (which is important, because it is much more likely you will sell the product if it is still in the user's cart when they return). However, if a user is routed to one node for a session, and then a different node on their next visit, there can be inconsistencies since the new node may be missing that user's cart contents. (Wouldn't you be upset if you put a 6 pack of Mountain Dew in your cart and then came back and it was empty?) One way around this can be to make sessions sticky so that the user is always routed to the same node, but then it is very hard to take advantage of some reliability features like automatic failover. In this case, the user's shopping cart would always have the contents, but if their sticky node became unavailable there would need to be a special case and the assumption of the contents being there would no longer be valid (although hopefully this assumption

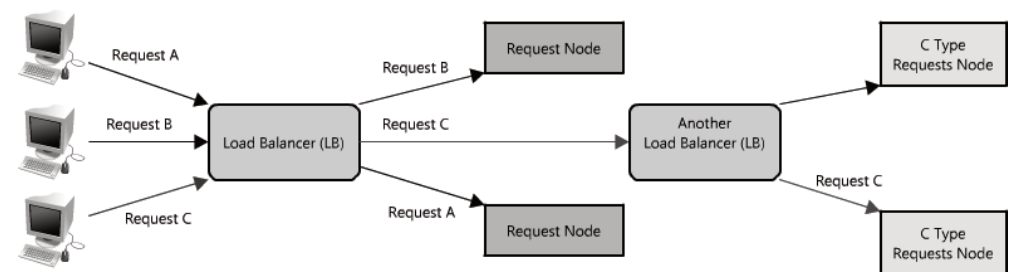


Figure 19: Multiple load balancers.

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

wouldn't be built into the application). Of course, this problem can be solved using other strategies and tools in this article, like services, and many not covered (like browser caches, cookies, and URL rewriting).

If a system only has a couple of a nodes, systems like round robin DNS may make more sense since load balancers can be expensive and

“In the cases where writes, or any task for that matter, may take a long time, achieving performance and availability requires building asynchrony into the system”

add an unneeded layer of complexity. Of course in larger systems there are all sorts of different scheduling and load-balancing algorithms, including simple ones like random choice or round robin, and more sophisticated mechanisms that take things like utilization and capacity into consideration. All of these algorithms allow traffic and requests to be distributed, and can provide helpful reliability tools like automatic failover, or automatic removal of a bad node (such as when it becomes unresponsive). However, these advanced features can make problem diagnosis cumbersome. For example, when it comes to high load situations, load balancers will remove nodes that may be slow or timing out (because of too many requests), but that only exacerbates the situation for the other nodes. In these cases extensive monitoring is important, because overall system traffic and throughput may look like

it is decreasing (since the nodes are serving less requests) but the individual nodes are becoming maxed out.

Load balancers are an easy way to allow you to expand system capacity, and like the other techniques in this article, play an essential role in distributed system architecture. Load balancers also provide the critical function of being able to test the health of a node, such that if a node is unresponsive or over-loaded, it can be removed from the pool handling requests, taking advantage of the redundancy of different nodes in your system.

Queues

So far we have covered a lot of ways to read data quickly, but another important part of scaling the data layer is effective management of writes. When systems are simple, with minimal processing loads and small databases, writes can be predictably fast; however, in more complex systems writes can take an almost non-deterministically long time. For example, data may have to be written several places on different servers or indexes, or the system could just be under high load. In the cases where writes, or any task for that matter, may take a long time, achieving performance and availability requires building asynchrony into the system; a common way to do that is with queues.

Imagine a system where each client is requesting a task to be remotely serviced. Each of these clients sends their request to the server, where the server completes the tasks as quickly as possible and returns the results to their respective clients. In small systems where one server (or logical service) can service incoming clients just as fast as they come, this sort of situation should work just fine. However, when the server receives more requests than it can handle, then each client is forced to wait for the other clients' requests to complete before a re-

IN THIS ISSUE

- [Editorial >>](#)
- [ASP.NET >>](#)
- [TypeScript >>](#)
- [Scalable Web >>](#)
- [Links >>](#)
- [Table of Contents >>](#)

sponse can be generated. This is an example of a synchronous request, depicted in Figure 20.

This kind of synchronous behavior can severely degrade client performance; the client is forced to wait, effectively performing zero work, until its request can be answered. Adding additional servers to address system load does not solve the problem either; even with effective load balancing in place it is extremely difficult to ensure the even and fair distribution of work required to maximize client performance. Further, if the server handling requests is unavailable, or fails, then the

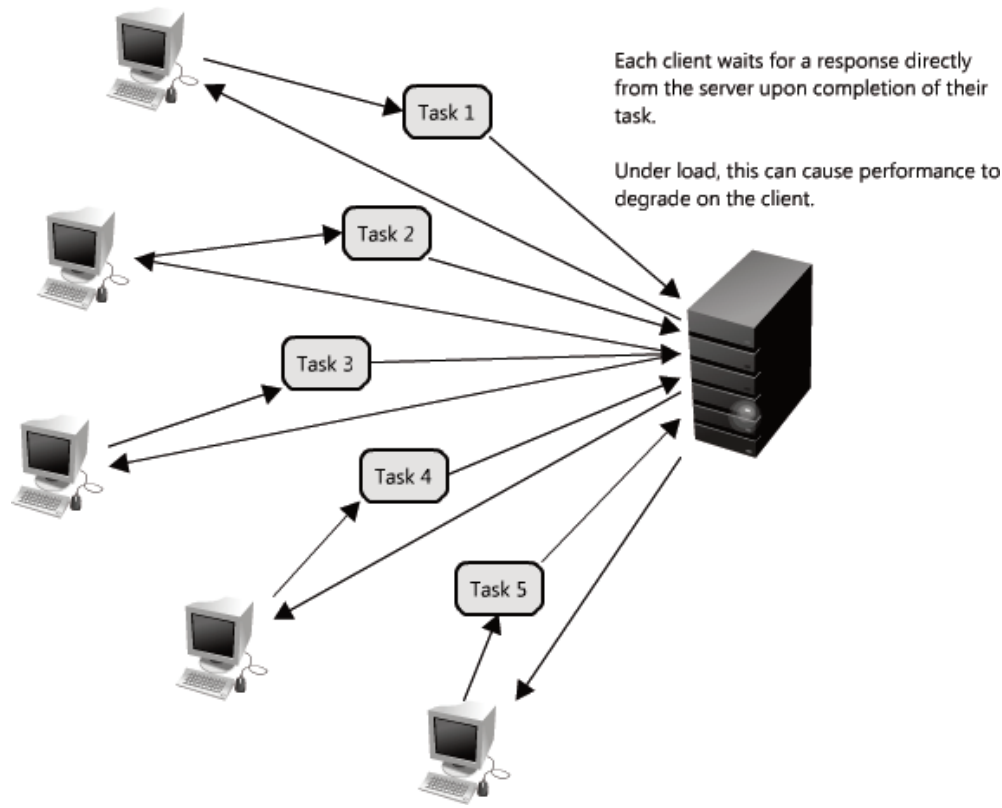


Figure 20: Synchronous request.

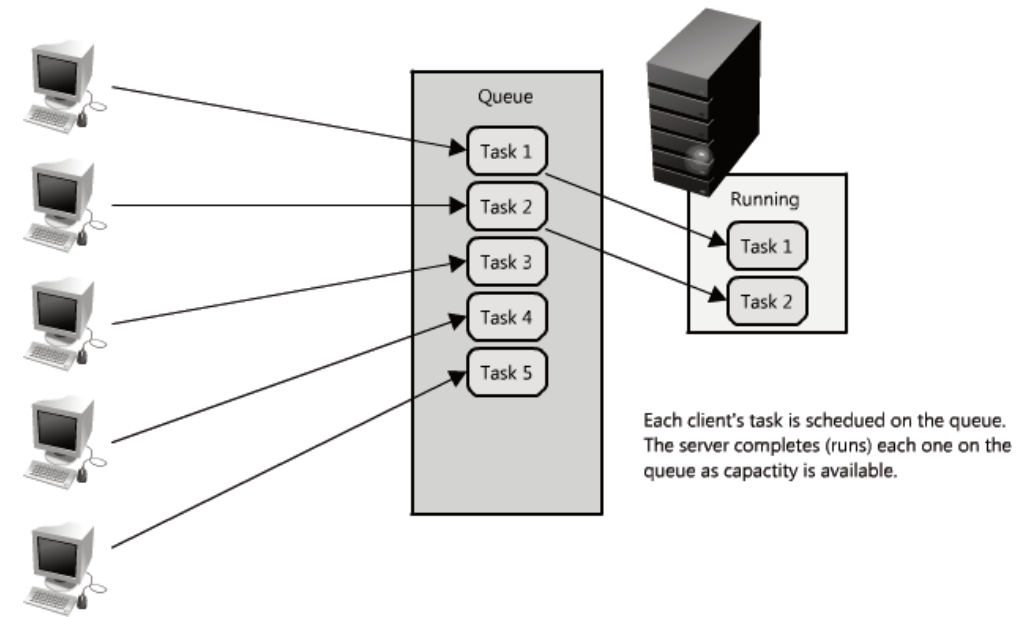


Figure 21: Using queues to manage requests.

clients upstream will also fail. Solving this problem effectively requires abstraction between the client's request and the actual work performed to service it.

Enter queues. A queue is as simple as it sounds: a task comes in, is added to the queue and then workers pick up the next task as they have the capacity to process it. (See Figure 21.) These tasks could represent simple writes to a database, or something as complex as generating a thumbnail preview image for a document. When a client submits task requests to a queue they are no longer forced to wait for the results; instead they need only acknowledgement that the request was properly received. This acknowledgement can later serve as a reference for the results of the work when the client requires it.

IN THIS ISSUE[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

Queues enable clients to work in an asynchronous manner, providing a strategic abstraction of a client's request and its response. On the other hand, in a synchronous system, there is no differentiation between request and reply, and they therefore cannot be managed separately. In an asynchronous system the client requests a task, the service responds with a message acknowledging the task was received,

“Designing efficient systems with fast access to lots of data is exciting, and there are lots of great tools that enable all kinds of new applications”

and then the client can periodically check the status of the task, only requesting the result once it has completed. While the client is waiting for an asynchronous request to be completed it is free to perform other work, even making asynchronous requests of other services. The latter is an example of how queues and messages are leveraged in distributed systems.

Queues also provide some protection from service outages and failures. For instance, it is quite easy to create a highly robust queue that can retry service requests that have failed due to transient server failures. It is more preferable to use a queue to enforce quality-of-service guarantees than to expose clients directly to intermittent service outages, requiring complicated and often-inconsistent client-side error handling.

Queues are fundamental in managing distributed communication between different parts of any large-scale distributed system, and there are lots of ways to implement them. There are quite a few open source queues such as RabbitMQ (<http://www.rabbitmq.com/>), ActiveMQ (<http://activemq.apache.org/>), and BeanstalkD (<http://is.gd/D7TdrA>), but some also use services like Zookeeper (<http://zookeeper.apache.org/>), or even data stores like Redis (<http://redis.io/>).

Conclusion

Designing efficient systems with fast access to lots of data is exciting, and there are lots of great tools that enable all kinds of new applications. This article covered just a few examples, barely scratching the surface, but there are many more — and there will only continue to be more innovation in the space.

— *Kate Matsudaira has worked as the VP Engineering/CTO at several technology startups, including currently at Decide, and formerly at SEOmoz and Delve Networks (acquired by Limelight). Prior to joining the startup world she spent time as a software engineer and technical lead/manager at Amazon and Microsoft. You can read more on her blog and website <http://katemats.com>.*

This work is made available under the Creative Commons Attribution 3.0 Unported license. Please see the full description of the license for details at <http://creativecommons.org/licenses/by/3.0/legalcode>. This article was excerpted with permission from *The Architecture of Open Source Applications, Volume II* (<http://is.gd/s9Gv1a>).

[Comment](#)

IN THIS ISSUE

[Editorial >>](#)[ASP.NET >>](#)[TypeScript >>](#)[Scalable Web >>](#)[Links >>](#)[Table of Contents >>](#)

This Month on DrDobbs.com

Items of special interest posted on www.drdobbs.com over the past month that you may have missed

AFTER XML, JSON: THEN WHAT?

Data representation, a seemingly simple problem, is proving intractable as one standard after another fails to deliver what is needed.

<http://www.drdobbs.com/240151851>

ADA 2012: ADA WITH CONTRACTS

The most important new feature in Ada 2012 is support for contract-based programming, which adds more validation of mission-critical code to a language already famous for its focus on reliability.

<http://www.drdobbs.com/240150569>

MONGODB WITH C#: DEEP DIVE

Using MongoDB in an ASP.NET MVC Web app.

<http://www.drdobbs.com/240152181>

THE QUIET REVOLUTION IN PROGRAMMING

During the last two years, one of the longest eras in programming has quietly drawn to a close.

<http://www.drdobbs.com/240152206>

NETREXX: THE ORIGINAL JVM SCRIPTING LANGUAGE RETURNS

The first scripting language for the JVM was recently open-sourced by IBM. After a major update, it is once again finding its place by leveraging its unique benefits.

<http://www.drdobbs.com/240149819>

NOSQL OPTIONS COMPARED

As NoSQL takes root in the enterprise, it's becoming clear that developers will need to gain greater familiarity with both the concepts behind it, as well as the various products that support it.

<http://www.drdobbs.com/240151198>

BUILDING AN AUTOMOTIVE DASHBOARD ON ANDROID

If you are an Android developer who spends a good deal of time in your car commuting and you want to get in on a rapidly expanding market for your Android coding skills, I highly recommend the Asteroid Smart as a great place to start.

<http://www.drdobbs.com/240152314>

IN THIS ISSUE

- [Editorial >>](#)
- [ASP.NET >>](#)
- [TypeScript >>](#)
- [Scalable Web >>](#)
- [Links >>](#)
- [Table of Contents >>](#)

Dr. Dobb's

Andrew Binstock Editor in Chief, Dr. Dobb's
andrew.binstock@ubm.com

Deirdre Blake Managing Editor, Dr. Dobb's
deirdre.blake@ubm.com

Amy Stephens Copyeditor, Dr. Dobb's
amy.stephens@ubm.com

Sean Coady Webmaster, Dr. Dobb's
sean.coady@ubm.com

Jon Erickson Editor in Chief Emeritus, Dr. Dobb's

CONTRIBUTING EDITORS

Scott Ambler
Mike Riley
Herb Sutter

DR. DOBB'S EDITORIAL
 751 Laurel Street #614
 San Carlos, CA
 94070
 USA

UBM TECH
 303 Second Street,
 Suite 900, South Tower
 San Francisco, CA 94107
 1-415-947-6000

INFORMATIONWEEK

Rob Preston VP and Editor In Chief, InformationWeek
rob.preston@ubm.com 516-562-5692

John Foley Editor, InformationWeek
john.foley@ubm.com 516-562-7189

Chris Murphy Editor, InformationWeek
chris.murphy@ubm.com 414-906-5331

Alexander Wolfe Editor In Chief, InformationWeek.com
alexander.wolfe@ubm.com 516-562-7821

Stacey Peterson Executive Editor, Quality, InformationWeek
stacey.peterson@ubm.com 516-562-5933

Lorna Garey Executive Editor, Analytics, InformationWeek
lorna.garey@ubm.com 978-694-1681

Stephanie Stahl Executive Editor, InformationWeek
stephanie.stahl@ubm.com 703-266-6030

Fritz Nelson VP and Editorial Director
fritz.nelson@ubm.com 949-223-3608

David Berlind Chief Content Officer, UBM Tech
david.berlind@ubm.com 978-462-5315

ART/DESIGN

Mary Ellen Forte Senior Art Director
maryellen.forte@ubm.com

Sek Leung Senior Designer
sek.leung@ubm.com

INFORMATIONWEEK.COM

Benjamin Tomkins Managing Editor
benjamin.tomkins@ubm.com 516-562-5336

Roma Nowak Senior Director, Online Operations and Production
roma.nowak@ubm.com 516-562-5274

Tom LaSusa Managing Editor, Newsletters

tom.lasusa@ubm.com

Jeanette Hafke Web Production Manager
jeanette.hafke@ubm.com

Joy Culbertson Web Producer
joy.culbertson@ubm.com

Nevin Berger Senior Director, User Experience
nevin.berger@ubm.com

Atif Malik Director, Web Development
atif.malik@ubm.com

INFORMATIONWEEK BUSINESS TECHNOLOGY NETWORK

EVP of Group Sales, InformationWeek Business Technology Network, Martha Schwartz
 (212) 600-3015, martha.schwartz@ubm.com

Sales Assistant, Salvatore Silletti
 (212) 600-3327, salvatore.silletti@ubm.com

SALES CONTACTS—WEST
 Western U.S. (Pacific and Mountain states) and Western Canada (British Columbia, Alberta)

Sales Director, Michele Hurabiell
 (415) 378-3540, michele.hurabiell@ubm.com

Strategic Accounts

Account Director, Sandra Kupiec
 (415) 947-6922, sandra.kupiec@ubm.com

Account Manager, Vesna Beso
 (415) 947-6104, vesna.beso@ubm.com

Account Executive, Matthew Cohen-Meyer
 (415) 947-6214, matthew.meyer@ubm.com

MARKETING

VP, Marketing, Winnie Ng-Schuchman
 (631) 406-6507, winnie.ng@ubm.com

Marketing Director, Angela Lee-Moll
 (516) 562-5803, angele.leemoll@ubm.com

Marketing Manager, Monique Kakegawa
 (949) 223-3609, monique.luttrell@ubm.com

Program Manager, Diane Scala
 516-562-5476, diane.scala@ubm.com

SALES CONTACTS—EAST

Midwest, South, Northeast U.S. and Eastern Canada (Saskatchewan, Ontario, Quebec, New Brunswick)

District Manager, Steven Sorhaindo
 (212) 600-3092, steven.sorhaindo@ubm.com

Strategic Accounts

District Manager, Mary Hyland
 (516) 562-5120, mary.hyland@ubm.com

Account Manager, Tara Bradeen
 (212) 600-3387, tara.bradeen@ubm.com

Account Manager, Jennifer Gambino
 (516) 562-5651, jennifer.gambino@ubm.com

Account Manager, Elyse Cowen
 (212) 600-3051, elyse.cowen@ubm.com

Sales Assistant, Kathleen Jurina
 (212) 600-3170, kathleen.jurina@ubm.com

AUDIENCE DEVELOPMENT

Director, Karen McAleer
 (516) 562-7833, karen.mcaleer@ubm.com

BUSINESS OFFICE

General Manager, Marian Dujmovits
United Business Media LLC
 600 Community Drive
 Manhasset, N.Y. 11030
 (516) 562-5000

Copyright 2013.
 All rights reserved.

UBM TECH

Paul Miller, CEO
Kathy Astromoff, CEO, Electronics
Robert Faletta, CEO, Channel
Kelley Damore, Chief Community Officer
Marco Pardi, President, Business Technology Events
David Berlind, Chief Content Officer
John Dennehy, Chief Financial Officer
David Michael, Chief Information Officer
Martha Schwartz, Chief Sales Officer, Business Technology Media
Scott Vaughan, Chief Marketing Officer
Simon Carless, EVP, Game & App Development and Black Hat
Lenny Heymann, EVP, New Markets
Angela Scalpello, SVP, People & Culture

UNITED BUSINESS MEDIA LLC

Pat Nohilly Sr. VP, Strategic Development and Business Administration
Marie Myers Sr. VP, Manufacturing

INFORMATIONWEEK VIDEO

informationweek.com/tv
Fritz Nelson Executive Producer
fritz.nelson@ubm.com

INFORMATIONWEEK BUSINESS TECHNOLOGY NETWORK

DarkReading.com
 Security
Tim Wilson, Site Editor
tim.wilson@ubm.com
IntelligentEnterprise.com
 App Architecture
Doug Henschen, Editor in Chief
doug.henschen@ubm.com

NetworkComputing.com

Networking, Communications, and Storage
Andrew Conry-Murray, Editor
andrew.murray@ubm.com
PlugIntoTheCloud.com
 Cloud Computing
John Foley, Site Editor
john.foley@ubm.com
InformationWeek SMB
 Technology for Small and Midsize Business
Benjamin Tomkins, Site Editor
benjamin.tomkins@ubm.com
Byte.com
Larry Seltzer
 Editorial Director
larry.seltzer@ubm.com
Dr. Dobb's
 Good Stuff for Serious Developers
Andrew Binstock
 Editor in Chief
andrew.binstock@ubm.com

Entire contents Copyright © 2013, UBM Tech/United Business Media LLC, except where otherwise noted. No portion of this publication may be reproduced, stored, transmitted in any form, including computer retrieval, without written permission from the publisher. All Rights Reserved. Articles express the opinion of the author and are not necessarily the opinion of the publisher. Published by UBM Tech/United Business Media, 303 Second Street, Suite 900 South Tower, San Francisco, CA 94107 USA 415-947-6000.