

moz://a

Diving into PDF.js

The advantages of the web platform

Calixte Denizet <calixte@mozilla.com>

Marco Castelluccio <marco@mozilla.com>

FOSDEM 2024

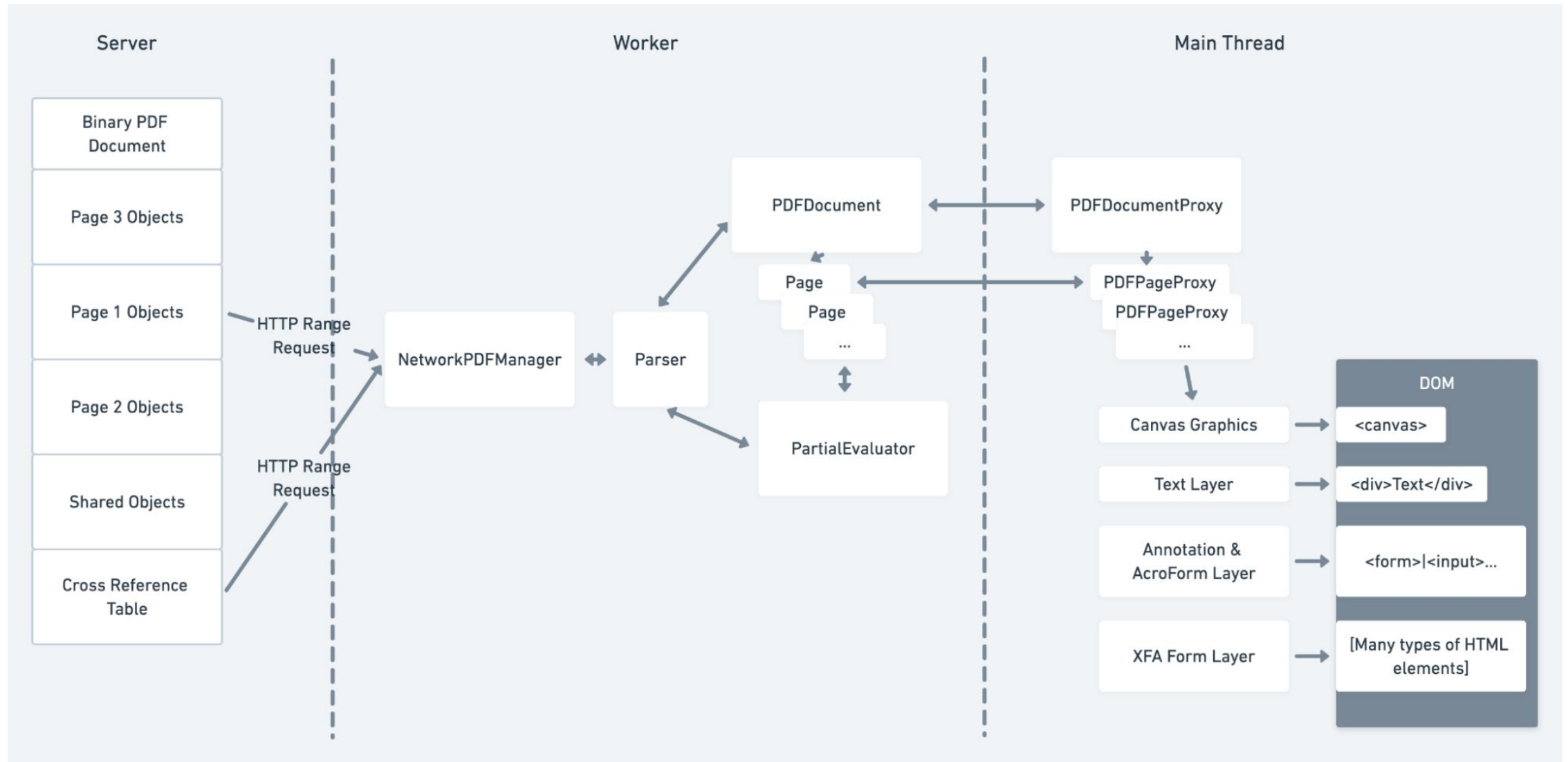
PDF.js

PDF format

- A PDF is a sequence of objects identified by a number and a table mapping identifiers and object offsets in the file
- Objects can be dictionaries, arrays, binary streams, ...
- PDF files can embed images, fonts, drawing instructions, xml, ...

PDF.js

Global architecture



How does it work?

- Data streamed into a worker
- The worker has a lot of parsers:
 - parser for the **pdf**
 - parsers for the embedded **images** (jpeg, jpeg2000, jbig2, ...)
 - parsers for the different **font** formats
 - **xml** parser
- Objects are parsed and sanitized (and fixed!)
- A drawing instructions list is generated and sent back to the main thread to be drawn on a HTML5 canvas.

Et voila !

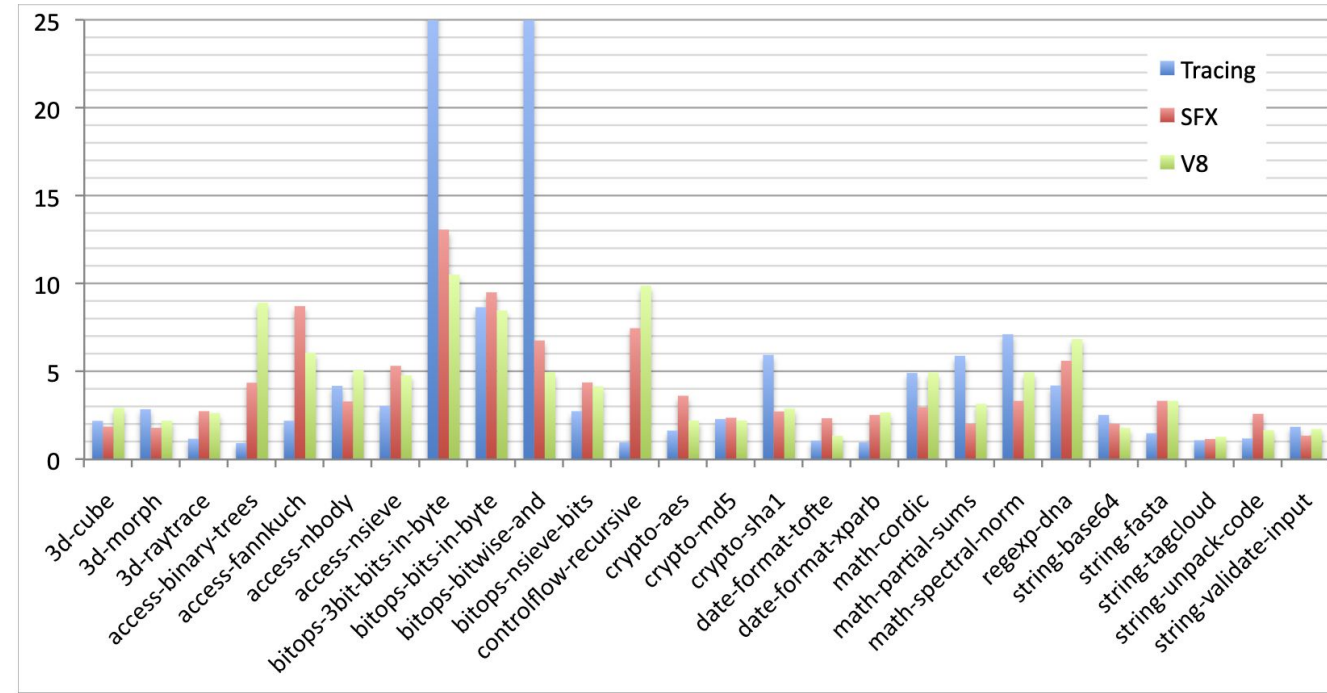


Figure 10. Speedup vs. a baseline JavaScript interpreter (SpiderMonkey) for our trace-based JIT compiler, Apple’s SquirrelFish Extreme inline threading interpreter and Google’s V8 JS compiler. Our system generates particularly efficient code for programs that benefit most from type specialization, which includes SunSpider Benchmark programs that perform bit manipulation. We type-specialize the code in question to use integer arithmetic, which substantially improves performance. For one of the benchmark programs we execute 25 times faster than the SpiderMonkey interpreter, and almost 5 times faster than V8 and SFX. For a large number of benchmarks all three VMs produce similar results. We perform worst on benchmark programs that we do not trace and instead fall back onto the interpreter. This includes the recursive benchmarks `access-binary-trees` and `control-flow-recursive`, for which we currently don’t generate any native code.

In particular, the `bitops` benchmarks are short programs that perform many bitwise operations, so TraceMonkey can cover the entire program with 1 or 2 traces that operate on integers. TraceMonkey runs all the other programs in this set almost entirely as native code.

- Two programs trace well, but have a long compilation time. `access-nbody` forms a large number of traces (81). `crypto-md5` forms one very long trace. We expect to improve performance on this programs by improving the compilation speed of nano-
...

Text layer

- To make the text selectable, there is a **text layer** with all the text

Trace-based Just-in-Time Type Specialization for Dynamic Languages

Andreas Gal^{*+}, Brendan Eich^{*}, Mike Shaver^{*}, David Anderson^{*}, David Mandelin^{*},
Mohammad R. Haghighat^{\$}, Blake Kaplan^{*}, Graydon Hoare^{*}, Boris Zbarsky^{*}, Jason Orendorff^{*},
Jesse Ruderman^{*}, Edwin Smith[#], Rick Reitmaier[#], Michael Bebenita⁺, Mason Chang^{+ #}, Michael Franz⁺

Mozilla Corporation^{*}

{gal, brendan, shaver, danderson, dmandelin, mrbkap, graydon, bz, jorendorff, jruderman}@mozilla.com

Adobe Corporation[#]

Annotation layer

- To render forms, an annotation layer is added with HTML form elements

		<input type="text"/>	<input type="text"/>
		code postal	localité
<input type="text"/>	au	<input type="text"/>	classe d'ir
tions brutes ²⁾		<input type="text"/>	12,00
Nature ³⁾		<input type="text"/>	13,00
		<input type="text"/>	
		<input type="text"/>	
	sous-total:	<input type="text"/>	25,00
s			n° dossier

Editing layer

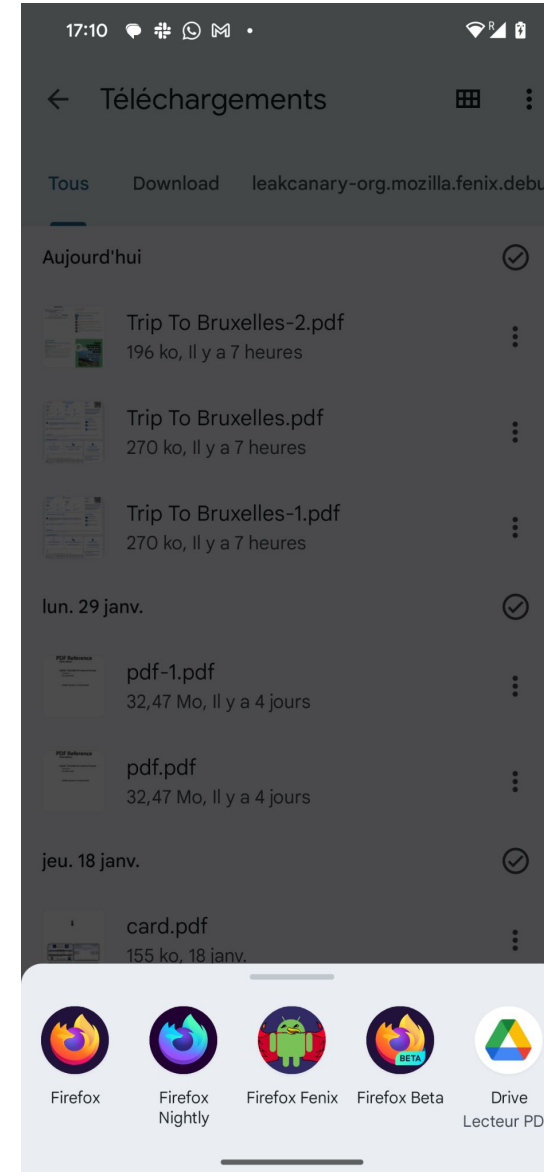
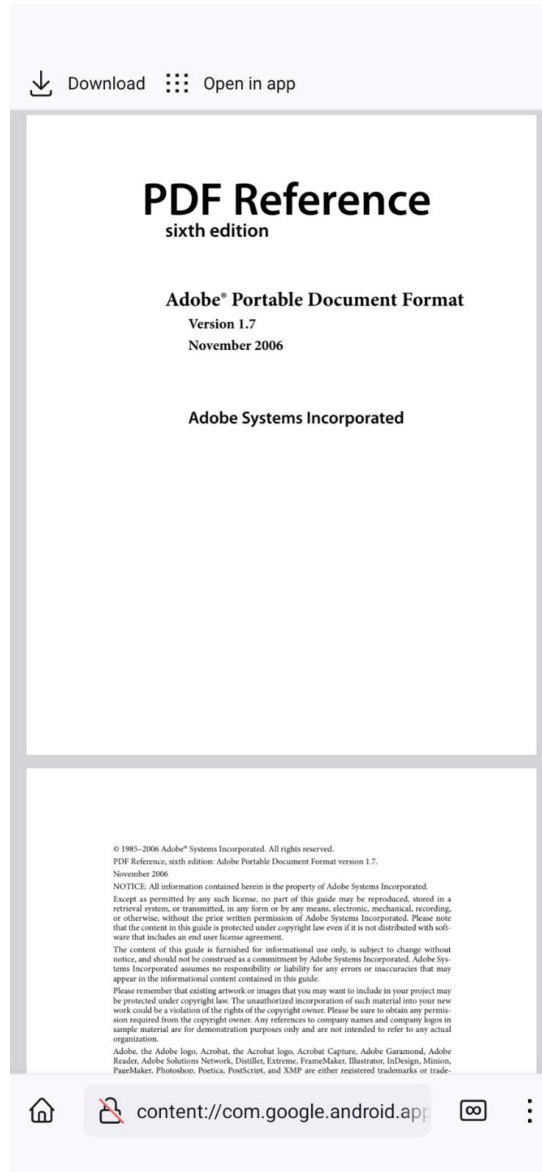
- To render added annotations, another layer is added



PDFs can embed JS

- To validate or compute form data, JS code can be embedded in PDFs.
- For the Firefox viewer, JS code is executed within a **sandbox**, i.e. a special isolated environment which has been used for years in WebExtensions.
- For the non-Firefox viewer, JS code is executed with **QuickJS** (a JS engine written in C) compiled in WASM.

Android



What we are doing

- Editing features

The screenshot displays a rich text editor with several elements:

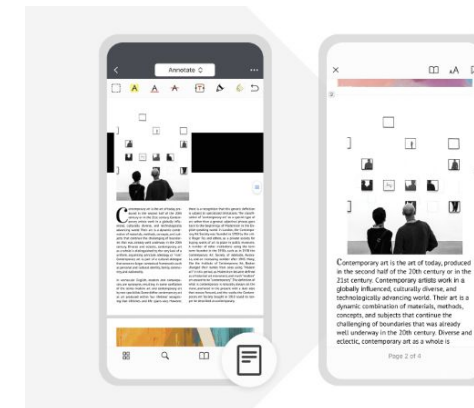
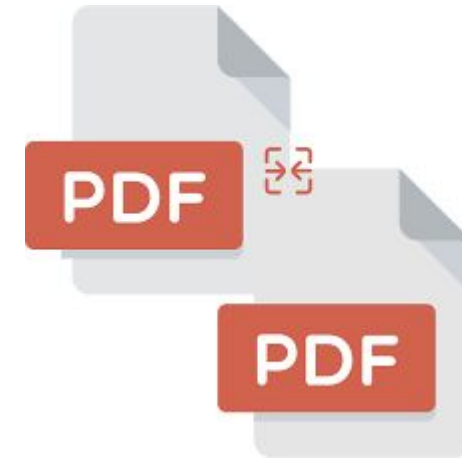
- Text:** "Hello World" in green font, enclosed in a blue border.
- Title:** "Trace-based Just-in-Time Type Specialization for Dynamic Languages" with "Type Specialization" highlighted in a light blue box.
- Authors:** Andreas Gal⁺, Brendan Eich*, Mike Shaver*, David Anderson*, David Mandelin*, Mohammad R. Haghighat^{\$}, Blake Kaplan*, Graydon Hoare*, Boris Zbarsky*, Jason Orendorff*, Jesse Ruderman*, Edwin Smith[#], Rick Reitmaier[#], Michael Bebenita⁺, Mason Chang^{+#}, Michael Franz⁺
- Emails:** Mozilla Corporation* {gal,brendan,shaver,danderson,dmandelin,mrbkap,graydon,bz,jorendorff,jruderman}@mozilla.com; Adobe Corporation[#] {edwsmith,rreitmai}@adobe.com; Intel Corporation {mohammad.r.haghighat}@intel.com; University of California, Irvine⁺ {mbebenit,changm,franz}@uci.edu
- Image:** Firefox logo with a blue border and a button labeled "+ Texte alternatif".
- Text:** "pdf.js" in blue cursive font, enclosed in a blue border.
- Panel:** A floating panel on the right with controls for "Couleur" (color), "Épaisseur" (thickness), and "Opacité" (opacity).
- Buttons:** Small trash icons are present next to the title and the "pdf.js" text.

Abstract
Dynamic languages such as JavaScript are more difficult to compile than statically typed ones. Since no concrete type information is available, traditional compilers need to emit generic code that can

and is used for the application logic of browser-based productivity applications such as Google Mail, Google Docs and Zimbra Collaboration Suite. In this domain, in order to provide a fluid user experience and enable a new generation of applications, virtual machines must provide a low startup time and high performance

What we are planning

- Merging and splitting, helping users not compromise their PDFs
- Signatures
- Reader mode



moz://a

Thank You