

Liquid IT : Toward a better compromise between development scalability and performance scalability not definitive

Etienne Brodu

November 23, 2015

Abstract

TODO translate from below when ready

Résumé

Internet étend nos moyens de communications, et réduit leur latence ce qui permet de développer l'économie à l'échelle planétaire. Il permet à chacun de mettre un service à disposition de milliards d'utilisateurs, en seulement quelques heures. La plupart des grands services actuels ont commencé comme de simples applications créées dans un garage par une poignée de personnes. C'est cette facilité à l'entrée qui a permis jusqu'à maintenant une telle croissance sur le web. Google, Facebook ou Twitter en sont quelques exemples. Au cours du développement d'une application, il est important de suivre cette croissance, au risque de se faire rattraper par la concurrence. Le développement est guidé par les besoins en terme de fonctionnalités, afin de vérifier rapidement si le service peut satisfaire l'audience. On parle d'approche modulaire des fonctionnalités. Des langages tel que Ruby ou Java se sont imposés comme les langages du web, justement parce qu'ils suivent cette approche qui permet d'intégrer facilement de nouvelles fonctionnalités.

Si une application répond correctement aux besoins, elle atteindra de manière virale un nombre important d'utilisateurs. Son audience peut prendre plusieurs ordres de grandeurs en quelques jours seulement, ou même en quelques heures suivant comment elle est relayée. Une application est dite *scalable* si elle peut absorber ces augmentations d'audience. Or il est difficile pour une application suivant l'approche modulaire d'être *scalable*.

Au moment où l'audience commence à devenir trop importante, il est nécessaire de modifier l'approche de développement de l'application. Le plus souvent cela implique de la réécrire complètement en utilisant des infrastructures *scalables* qui imposent des modèles de programmation et des API spécifiques, qui représentent une charge de travail conséquente et incertaine. De plus, l'équipe de développement doit concilier cette nouvelle approche de développement *scalable*, avec la demande en fonctionnalités. Aucun langage n'a clairement réussi le compromis entre ces deux objectifs.

Pour ces raisons, ce changement est un risque pour la pérennité de l'application. D'autant plus que le cadre économique accorde peu de marges d'erreurs, comme c'est le cas dans la plupart des start-up, mais également dans de plus grandes structures.

Cette thèse est source de propositions pour écarter ce risque. Elle repose sur les deux observations suivantes. D'une part, Javascript est un langage qui a gagné en popularité ces dernières années. Il est omniprésent sur les clients, et commence à s'imposer également sur les serveurs avec Node.js. Il a accumulé une communauté de

développeurs importante, et constitue l'environnement d'exécution le plus largement déployé. De ce fait, il se place maintenant de plus en plus comme le langage principal du web, détrônant Ruby ou Java. D'autre part, l'exécution de Javascript s'assimile à un pipeline. La boucle événementielle de Javascript exécute une suite de fonctions dont l'exécution est indépendante, mais qui s'exécutent sur un seul cœur pour profiter d'une mémoire globale.

L'objectif de cette thèse est de maintenir une double représentation d'un code Javascript grâce à une équivalence entre l'approche modulaire, et l'approche pipeline d'un même programme. La première répondant aux besoins en fonctionnalités, et favorisant les bonnes pratiques de développement pour une meilleure maintenabilité. La seconde proposant une exécution plus efficace que la première en permettant de rendre certaines parties du code relocalisables en cours d'exécution.

Nous étudions la possibilité pour cette équivalence de transformer un code d'une approche vers l'autre. Grâce à cette transition, l'équipe de développement peut continuellement itérer le développement de l'application en suivant les deux approches à la fois, sans être cloisonné dans une, et coupé de l'autre.

Nous construisons un compilateur permettant d'identifier les fonctions de Javascript et de les isoler dans ce que nous appelons des Fluxions, contraction entre fonctions et flux. Un conteneur qui peut exécuter une fonction à la réception d'un message, et envoyer des messages pour continuer le flux vers d'autres fluxions. Les fluxions sont indépendantes, elles peuvent être déplacées d'une machine à l'autre.

Nous montrons qu'il existe une correspondance entre le programme initial, purement fonctionnel, et le programme pivot fluxionnel afin de maintenir deux versions équivalentes du code source. En ajoutant à un programme écrit en Javascript son expression en Fluxions, l'équipe de développement peut le rendre *scalable* sans effort, tout en étant capable de répondre à la demande en fonctionnalités.

Ce travail s'est fait dans le cadre d'une thèse CIFRE dans la société Worldline. L'objectif pour Worldline est de se maintenir à la pointe dans le domaine du développement et de l'hébergement logiciel à travers une activité de recherche. L'objectif pour l'équipe Dice est de conduire une activité de recherche en partenariat avec un acteur industriel.

Contents

1	Introduction	2
1.1	Web development	3
1.2	Performance requirements	3
1.3	Problematic and proposal	4
1.4	Thesis organization	5
2	Context and objectives	6
2.1	The Web as a Platform	7
2.1.1	Javascript, The Language of the Web	7
2.1.1.1	Historical Context	8
2.1.1.2	Current Situation	10
2.1.1.3	Event-Loop Execution Model	14
2.1.2	Highly Concurrent Web Servers	16
2.1.2.1	Scalable Concurrency	16
2.1.2.2	Time-slicing and Parallelism	17
2.1.2.3	Pipeline Execution Model	17
2.2	An Economical Problem	18
2.2.1	Disrupted Development	18
2.2.1.1	Power-Wall Disruption	18
2.2.1.2	Unavoidable Modularity	19
2.2.1.3	Technological Shift	19
2.2.2	Seamless Web Development	20
2.2.2.1	Real-Time Streaming Web Services	20
2.2.2.2	Differences	20
2.2.2.3	Equivalence	21
3	Software Design, State Of The Art	22
3.1	Software Maintainability	24

3.1.1	Modular Programming	25
3.1.1.1	Design Choices	25
3.1.1.2	Programming Models	26
3.1.1.3	Performance Limitations	28
3.1.2	Performance Improvements	29
3.1.2.1	Concurrent Programming	30
3.1.2.2	Compilation	32
3.1.2.3	Accessibility Limitations	35
3.2	Software Performance	35
3.2.1	Parallel Programming	36
3.2.1.1	Asynchronous and Isolated Process Parallelism	37
3.2.1.2	Stream Processing Systems	38
3.2.1.3	Elitism of Parallel Programming	39
3.2.2	Maintainability Improvements	40
3.2.2.1	Execution Organization	41
3.2.2.2	Memory Abstraction	42
3.2.2.3	Lack of Higher-Order Programming	42
3.3	Seamless Development	45
3.3.1	Equivalence	47
3.3.1.1	Rupture Point	47
3.3.1.2	Invariance	48
4	Pipeline extraction	49
4.1	Definitions	50
4.1.1	Callback	50
4.1.2	Promise	52
4.1.3	From continuations to Promises	53
4.1.4	Due	55
4.2	Equivalence	56
4.2.1	Execution order	56
4.2.2	Execution linearity	57
4.2.3	Variable scope	58
4.3	Compiler	58
4.3.1	Identification of continuations	58
4.3.2	Generation of chains	59
4.4	Evaluation	59

5 Pipeline isolation	62
5.1 Fluxional execution model	62
5.1.1 Fluxions	63
5.1.2 Messaging system	63
5.1.3 Service example	64
5.2 Fluxionnal compiler	67
5.2.1 Analyzer step	67
5.2.1.1 Rupture points	68
5.2.1.2 Detection	69
5.2.2 Pipeliner step	69
5.3 Real case test	71
5.3.1 Compilation	72
5.3.2 Isolation	73
5.3.2.1 Variable <code>req</code>	73
5.3.2.2 Closure <code>next</code>	73
5.3.3 Future works	75
6 Conclusion	76
A Language popularity	77
A.1 PopularitY of Programming Languages (PYPL)	77
A.2 TIOBE	78
A.3 Programming Language Popularity Chart	79
A.4 Black Duck Knowledge	79
A.5 Github	81
A.6 HackerNews Poll	81

Chapter 1

Introduction

Contents

1.1	Web development	3
1.2	Performance requirements	3
1.3	Problematic and proposal	4
1.4	Thesis organization	5

When the amazed 7 years old I was laid eyes on the first family computer, my life goal became to know everything there is to know about computers. This thesis is a mild achievement. It compiles my PhD work on *bridging the gap between development scalability, and performance scalability, in the case of real-time web applications*.

This work is the fruit of a collaboration between the Worldline and the Inria DICE (Data on the Internet at the Core of the Economy) team from the CITI (Centre d’Innovation en Télécommunications et Intégration de services) at INSA de Lyon. For Worldline, this work fall within a larger work named Liquid IT, on the future of the cloud infrastructure and development. As defined by Worldline, Liquid IT aims at decreasing the time to market of a web service, allows the development team to focus on service specifications rather than technical twists and ease service maintenance. The purpose of my work, was to separate development scalability from performance scalability, to allow a continuous development from prototyping phase, until runtime on thousands of clusters. On the other hand, the DICE team focuses on the consequences of technology on economical and social changes at the digital age. This work falls within this scope as the development of web services is driven by economical factors. Indeed, this thesis study the relation between the economical, and the technological constraints on the development of web services.*



This sentence
is to be dis-
cussed about.

1.1 Web development

The growth of web platforms is partially due to Internet’s capacity to allow very quick releases of a minimal viable product (MVP). In a matter of hours, it is possible to release a prototype and start gathering a user community around. “*Release early, release often*”, and “*Fail fast*” are the punchlines of the web entrepreneurial community. It is crucial for the prosperity of such project to quickly validate that the proposed solution meets the needs of its users. Indeed, the lack of market need is the first reason for startup failure.¹ That is why the development team quickly concretises an MVP and iterates on it using a feature-driven, monolithic approach. Such as proposed by imperative languages like Java or Ruby.

1.2 Performance requirements

If the service successfully complies with users requirements, its community might grow with its popularity. If it can quickly respond to this growth, it is scalable. However, it is difficult to develop scalable applications with the feature-driven approach

¹<https://www.cbinsights.com/blog/startup-failure-post-mortem/>

mentioned above. Eventually this growth requires to discard the initial monolithic approach to adopt a more efficient processing model instead. Many of the most efficient models distribute the system on a cluster of commodity machines.

Once split, the service parts are connected by an asynchronous messaging system. Many tools have been developed to express and manage these service parts and their communications. However, these tools impose specific interfaces and languages, different from the initial monolithic approach. It requires the development team either to be trained or to hire experts, and to start over the initial code base. This shift causes the development team to spend development resources in background without adding visible value for the users. It is a risk for the evolution of the project as the second and third reasons for startup failures are running out of cash, and missing the right competences.

1.3 Problematic and proposal

These shifts are a risk for the economical evolution of a web application by disrupting the continuity in its development. The main question I address in this thesis is how to avoid these shifts, so as to allow a continuous development? That is to reconcile the reactivity required in the early stage of development and the performance increasingly required with the growth of popularity. To answer this question, this thesis proposes a solution based on an equivalence between two different programming paradigms. On one hand, there is the imperative, functional, asynchronous programming model, embodied by Javascript. On the other hand, there is the dataflow, distributed, programming model, embodied by the concept of fluxions introduced in chapter 5.

This thesis contains two main contributions. The first contribution is a compiler allowing to split a program into a pipeline of stages depending on a common memory store. The second contribution, stemming from the first one, is a second compiler, allowing to make independent the stages of this pipeline. With these two contributions, it is possible to build a compiler that links an imperative representation with a flow-based representation. The imperative representation carries the functional modularization of the application, while the flow-based representation carries its execution distribution. A development team shall then use these two representations to continuously iterate over the implementation of an application, while keeping both maintainability and performance.

1.4 Thesis organization

This thesis is organized in four main chapters. Chapter 2 introduces the context for this thesis and explains in greater details its objectives. It presents the challenge to build web applications at a world wide scale, without jamming the organic evolution of its implementation. It concludes drawing a first answer to this challenge. Chapter 3 presents the works surrounding this thesis, and how they relate to it. It defines into the notions outlined in the precedent chapter to help the reader understand better the context. The end of this chapter presents clearly the problematic addressed in this thesis. Chapter 4 presents the first contribution allowing to represent a program as a pipeline of stages. It introduces Dues to encapsulate these stages, based on Javascript Promises. Chapter 5 presents the second contribution allowing to make these stages independent. It introduces Fluxion to encapsulate these stages. Chapter 6 concludes this thesis, and draws the possible perspectives beyond this work.

Chapter 2

Context and objectives

Contents

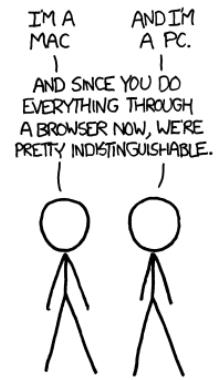
2.1	The Web as a Platform	7
2.1.1	Javascript, The Language of the Web	7
2.1.1.1	Historical Context	8
2.1.1.2	Current Situation	10
2.1.1.3	Event-Loop Execution Model	14
2.1.2	Highly Concurrent Web Servers	16
2.1.2.1	Scalable Concurrency	16
2.1.2.2	Time-slicing and Parallelism	17
2.1.2.3	Pipeline Execution Model	17
2.2	An Economical Problem	18
2.2.1	Disrupted Development	18
2.2.1.1	Power-Wall Disruption	18
2.2.1.2	Unavoidable Modularity	19
2.2.1.3	Technological Shift	19
2.2.2	Seamless Web Development	20
2.2.2.1	Real-Time Streaming Web Services	20
2.2.2.2	Differences	20
2.2.2.3	Equivalence	21

This chapter presents the general context for this work, and leads to a definition of the scope of this thesis. Section 2.1 presents the context of web development, and the motivations that led the web to become a software platform. It presents Javascript as the trending language currently taking over web development. Then, it presents the challenges of developing web servers for large audiences, after the prototyping phase. Section 2.2 states the problem tackled by this thesis, and its objectives. In the economical context of the Web, the languages often fail to grow with the project they initially supported very efficiently. The inadequacy of the languages to support the growth of web applications leads to wasted development efforts, and additional costs. The objective of this thesis is to avoid these efforts and costs. It intends to provide a continuous development from the initial prototype up to the releasing and maintenance of the complete product.

2.1 The Web as a Platform

Similarly to operating systems, Web browsers started as software products with extension capabilities with scripts and applications. The distribution of an applications is limited only by the platform it can be deployed on. The Web spreads the scalability of software distribution world wide with a near zero latency. It eventually became the main distribution medium for software, and the wider market there can possibly be. It led the Web to become the platform, replacing operating systems.

Now, with web services, or Software as a Service (SaaS), the distribution medium of software is so transparent that owning a software product to have an easier access is no longer relevant. It stimulates a completely new business model based on a free access for the user, while claiming value for their data. The next paragraphs present Javascript, the language that allowed this new business model to emerge.



2.1.1 Javascript, The Language of the Web

In the 80's with Moore's law predicting exponential increase in hardware performance, reducing development time became more profitable than reducing hardware costs. Higher-level languages replaced lower-level languages, because the economical gain in development time compensated the decrease in performance. Most of the now popular programming languages were released at this time, Python(1991), Ruby(1993), Java(1994), PHP(1995) and Javascript(1995).

Java thrived in the software industry. However, Java lost the hype that drove the community innovation and creativity, and now struggles to keep up with the latest trends in software development. On the contrary, Ruby on Rails emerged from an industrial context, but is now open source, and backed by a strong community that makes it evolve and mature. Other languages like Python and PHP, bloomed, grew within a strong community, and were later adopted by the industry for web development. Django, the Python web frameworks, is used to develop many web applications in industrial contexts. And Wordpress, a PHP publishing platform, is an economical success. These examples show that the involvement of the community is critical for the adoption, evolution and maturation of a language.

Since a few years, Javascript is slowly becoming the main language for web development. It is the only choice in the browser. This position became an incentive to make it fast (V8, ASM.js) and convenient (ES6, ES7). And since 2009, it is present on the server as well with Node.js. This omnipresence became an advantage. It allows to develop and maintain the whole application with the same language.

2.1.1.1 Historical Context

“There are only two kinds of languages: the ones people complain about and the ones nobody uses”

—B. Stroustrup¹

Javascript was created by Brendan Eich at Netscape around May 1995, and released to the public in September. At the time, Java was quickly adopted as the default language for web servers development, and everybody was betting on pushing Java inside the browser as well. The history proved them wrong².

Javascript was released as a scripting engine on Netscape Navigator and later on its concurrent, Internet Explorer. The competition between the two was fragmenting the Web. Web pages had to be designed for a specific browser. To stop this fragmentation, Netscape submitted Javascript to Ecma International for standardization in November 1996. ECMA International released ECMA-262 in June 1997, the first standard for Javascript - or ECMAScript. A standard to which all browser should refer for their implementations.

The initial release of Javascript was designed in a rush, within 10 days, and targeted unexperienced developers. For these reasons, the language was considered poorly designed and unattractive by the developer community.

illustration:
the ugly
duckling

²Except for mobiles, with the success of Dalvik Android. Though, firefox OS is emerging.

Why does Javascript suck?³ Is Javascript here to stay?⁴ Why Javascript Is Doomed.⁵ Why JavaScript Makes Bad Developers.⁶ JavaScript: The World's Most Misunderstood Programming Language⁷ Why Javascript Still Sucks⁸ 10 things we hate about JavaScript⁹ Why do so many people seem to hate Javascript?¹⁰

But things evolved drastically since. All web browsers include a Javascript interpreter, making Javascript the most ubiquitous runtime [26]. Any Javascript code is open, allowing the community to pick, improve and reproduce the best techniques¹¹. Javascript is distributed freely, with all the tools needed to reproduce and experiment on the largest communication network in history. All these reasons made the popularity of the Web and Javascript.

“When JavaScript was first introduced, I dismissed it as being not worth my attention. Much later, I took another look at it and discovered that hidden in the browser was an excellent programming language.”

—Douglas Crockford

Javascript was initially limited to short interactions on web pages. The typical usage was to pre-validate forms on the client to avoid wasting wrongly formated requests to the server. This situation hugely improved since the beginning of the language. Nowadays, there is a lot of web-based application replacing desktop applications, like mail client, word processor, music player, graphics editor...

ECMA International released several version in the few years following the creation of Javascript. The third version contributed to give Javascript a more complete and solid base as a programming language. From this point on, the consideration for Javascript kept improving.

In 2005, James Jesse Garrett released *Ajax: A New Approach to Web Applications*, a white paper coining the term Ajax [29]. It uses Javascript to dynamically

³<http://whydoesitsuck.com/why-does-javascript-suck/>

⁴<http://www.javaworld.com/article/2077224/learn-java/is-javascript-here-to-stay-.html>

⁵<http://simpleprogrammer.com/2013/05/06/why-javascript-is-doomed/>

⁶<https://thorprojects.com/blog/Lists/Posts/Post.aspx?ID=1646>

⁷<http://www.crockford.com/javascript/javascript.html>

⁸<http://www.boronine.com/2012/12/14/Why-JavaScript-Still-Sucks/>

⁹<http://www.infoworld.com/article/2606605/javascript/146732-10-things-we-hate-about-JavaScript.html>

¹⁰<https://www.quora.com/Why-do-so-many-people-seem-to-hate-JavaScript>

¹¹<http://blog.codinghorror.com/the-power-of-view-source/>

reload the content inside a web page, hence improving the user experience. It allows Javascript to develop richer applications inside the browser, from user interactions to network communications. The first web applications to use Ajax were Gmail, and Google maps¹². The community released Javascript framework to assist the development of these larger applications. Prototype¹³ and DOJO¹⁴ are early famous examples, and later jQuery¹⁵ and underscore¹⁶.

In 2004, the Web Hypertext Application Technology Working Group¹⁷ was formed to work on the fifth version of the HTML standard. The name is misleading, it is really about giving Javascript superpowers like geolocation, storage, audio, video, and many more. The releases of HTML5 and ECMAScript 5, in 2008 and 2009, represent a mile-stone in the development of web-based applications. Around the same time, Google released the Javascript interpreter V8 for its browser Chrome, improving drastically the execution performance. Javascript became the *defacto* programming language to develop on this rising application platform that is the Web.

2.1.1.2 Current Situation

The rise of Javascript is indisputable on the web, and seems to be rising in the software industry as well. But it is difficult to give an accurate representation of the situation because the software industry often maintains a fog of war to try to keep an edge. The following paragraphs report some efforts to clear up the situation.

Available Resources According to the TIOBE Programming Community index, Javascript ranks 8th, as of October 2015, and was the most rising language in 2014. This index measure the popularity of a programming language with the number of results on many search engines. However, this measure is controversial as the number of pages doesn't represent the number of readers. Alternatively, Javascript ranks 7th on the PYPL, as of October 2015. The PYPL index is based on Google trends to measure the number of requests on a programming language. However, it is limited to Google searches. From these indexes, the major programming languages are

¹²A more in-depth analysis of the history of Ajax, given by late Aaron Swartz <http://www.aaronsw.com/weblog/ajaxhistory>

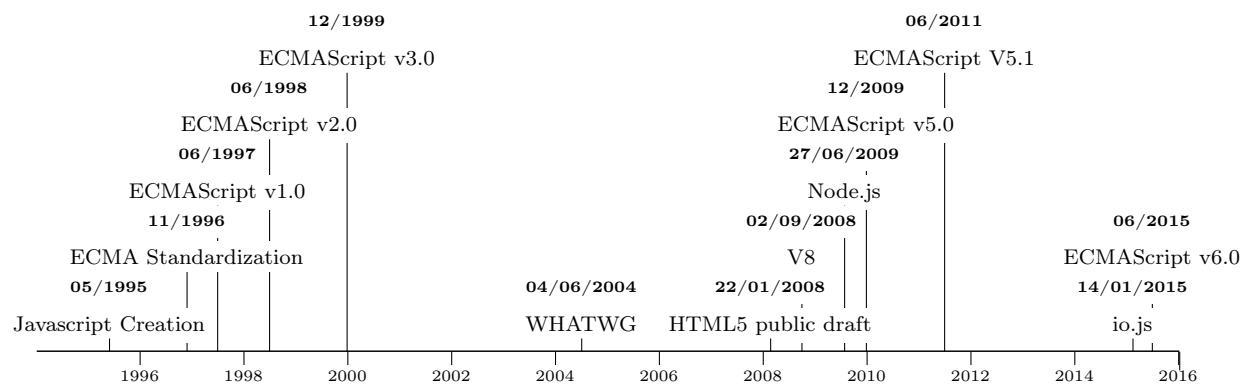
¹³<http://prototypejs.org/>

¹⁴<https://dojotoolkit.org/>

¹⁵<https://jquery.com/>

¹⁶<http://underscorejs.org/>

¹⁷<https://whatwg.org/>



Java, then C/C++, C# and Python. Javascript seems not as popular as previously described. The following paragraphs rectify this vision.

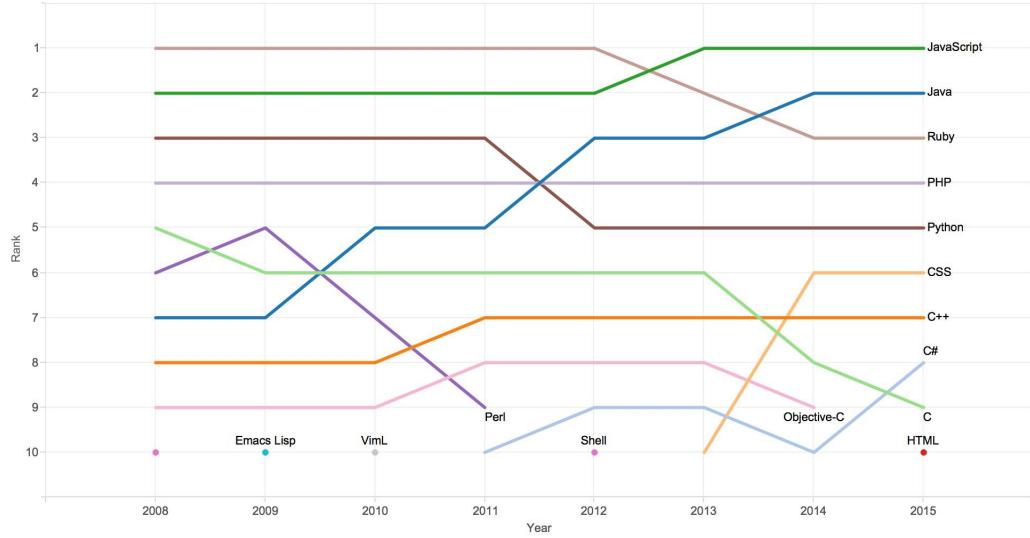
*



TODO graphical ranking of TIOBE and PYPL

Developers Collaboration Platforms Online collaboration tools gives an indicator of the number of developers and project using certain languages. Javascript is the most used language on *Github*, the most important collaborative development platform, with around 9 millions users. It represents more than 320 000 repositories, while the second language is Java with more than 220 000 repositories. Javascript is the most cited language on *StackOverflow*, the most important Q&A platform for developers. It represent more than 960 000 questions, while the second is Java with around 940 000 questions. Additionnaly, Javascript has the most important and impressive package repository growth. Moreover, Javascript is currently the second language used in open source projects, according to *Black Duck Software*¹⁸. C is first, C++ third and Java fourth.¹⁹ These four languages represent about 80% of all programming language usage in open source communities.

Rank of top languages on GitHub.com over time



Source: GitHub.com²⁰

*



TODO redo this graph, it is ugly.

*

¹⁸<https://www.blackducksoftware.com/>

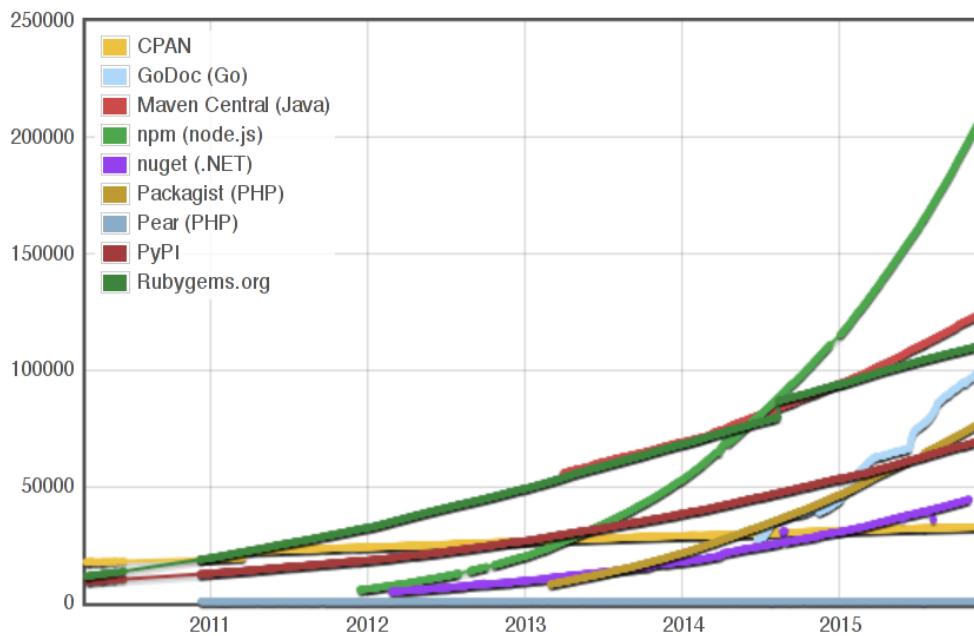
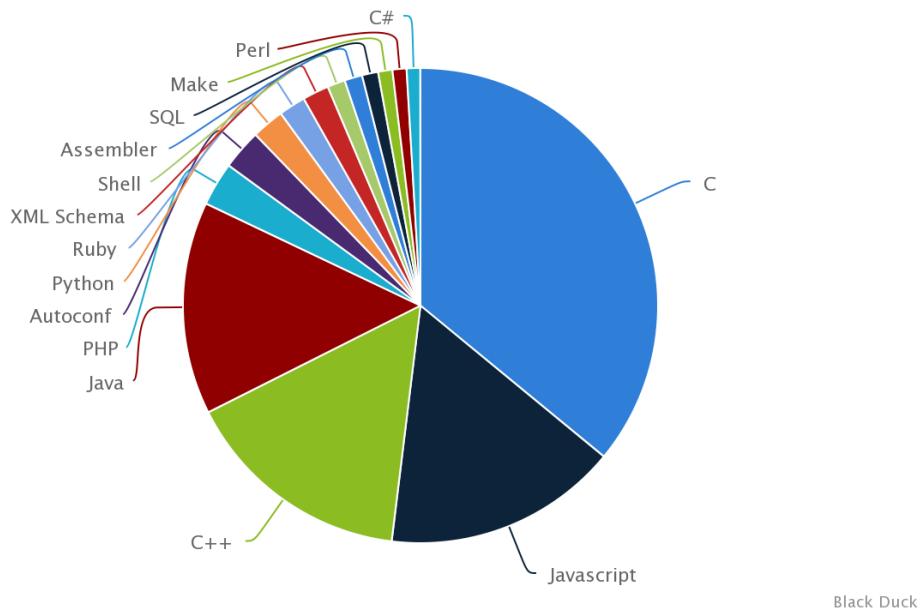
¹⁹<https://www.blackducksoftware.com/resources/data>

²⁰<https://github.com/blog/2047-language-trends-on-github>



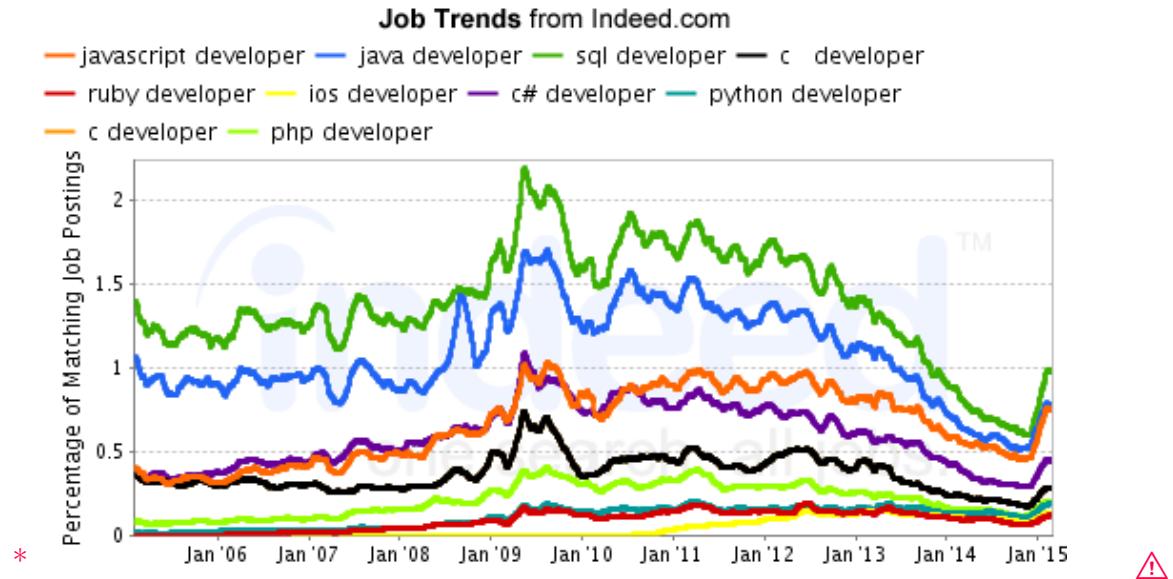
TODO graphical ranking of the tags in StackOverflow

Releases within the last 12 months



Jobs The actors of the software industry tends to hide their activities trying to keep an edge on the competition. The previous metrics represent the visible activity but

are barely representative of the software industry. The trends on job opportunities give some additional hints on the situation. Javascript is the third most wanted skill, according to *Indeed*²¹, right after SQL and Java.²² Moreover, according to *breaz.io*²³, Javascript developers get more opportunities than any other developers. Javascript is increasingly adopted in the software industry.



TODO redo
this graph, it
is ugly.

All these metrics represent different faces of the current situation of the Javascript adoption in the development community and industry. It is widely used on the web, in open source projects, and in the software industry. With the increasing importance of client web applications, Javascript is assuredly one of most important language in the times to come. Moreover, Javascript is also allows to build the server side of web applications. The next paragraphs presents the event-loop model used to develop Javascript web applications, both client and server-side.

2.1.1.3 Event-Loop Execution Model

Javascript is often associated with an event-based paradigm to react to concurrent user interactions. In 2009, Joyent released Node.js to build real-time web services with this paradigm. It is a server-side implementation of Javascript based on an

²¹<http://www.indeed.com>

²²<http://www.indeed.com/jobtrends?q=Javascript%2C+SQL%2C+Java%2C+C%2B%2B%2C+C%2FC%2B%2B%2C+C%23%2C+Python%2C+PHP%2C+Ruby&l=>

²³<https://breaz.io/>

event-loop. This event-based paradigm proved to be very efficient as well for a web service to react to concurrent requests. This section presents the event-loop execution model, and the advantages of Javascript for this paradigm.

The event-loop efficiency comes from non-blocking communications, asynchronous execution, and cooperative scheduling. It relies on a queue storing the messages received asynchronously. The loop executes previously defined tasks to process these messages one after the other. Each task can initiate new communications, leading in turn to the queuing of more messages, which trigger more tasks, and so on. Each task is executed atomically and exclusively, until it yields execution, to continue with the next task in queue.

*



TODO
schema of an
event-loop

Callbacks In Javascript, the asynchronous communications are initiated by function calls. This asynchronous callee immediately returns to avoid waiting the result. The task to process the result of the communication, and to continue the execution is a function passed as an argument to the callee. This function is name a callback or a continuation. A callback is a function passed as an argument to a callee, for the callee to transfer the control back to the caller after its execution, without the need for synchronization.

In this execution model, the execution is interrupted by asynchronous function calls. It organizes the execution of callbacks causally, one after the other. The input stream of data flows through a sequence of callbacks until the application outputs it. The asynchronous execution flow is controlled by callbacks, and is organized similarly to a pipeline. In this model, callbacks are the atoms of asynchronous execution flow control. The next paragraph presents a more elaborate form of control.

Promises Since the asynchronous execution flow became more complicated on larger web application, many projects proposed improved asynchronous execution controls on top of callbacks. The ECMAScript specification proposes Promises for such purpose. It arranges sequence of causally related callbacks into cleanly organized pipeline of callbacks communicating their result to the next.

Callbacks are said to be first class citizen. They imply higher-order programming, which is part of functional programming. Javascript features higher-order functions.

Closures For a callback to continue the execution without needing synchronization with the callee, it needs to have access to the initial context of the caller. This context is linked with the function when passed to the callee. The association of a function and its initial context is called a closure.

Higher-order programming is convenient for developers, as they allow great modularity in the implementation through *e.g.* inversion of control. It is presented in further details in section 3.1. However, because the contexts are passed, and shared all over the implementation, this programming model needs a global memory for coordination. As presented in the next section, this need is problematic to increase the concurrency of the execution.

This section presented Javascript as the language of the web, and its programming model. The next section presents the realities and technical challenges to assure the performance of web services against billions of users.

2.1.2 Highly Concurrent Web Servers

The previous section presented Javascript, the prolific language to build the Web. With SaaS, a Web service can scale world wide with near zero latency, and accessing it is as simple as distributing it world wide. With this broad range of distribution, a new business model emerged, allowing free access for the user. The usage exploded, and the software industry needed innovative solutions to cope with large network traffic.

2.1.2.1 Scalable Concurrency

The Internet allows communication at an unprecedented scale. There is more than 16 billions connected devices, and it is growing fast²⁴ [Hilbert2011]. A large web application like google search receives about 40 000 requests per seconds²⁵. Such a Web application needs to be highly concurrent to manage this amount of simultaneous requests. In the 2000s, the limit to break was 10 thousands simultaneous connections with a single commodity machine²⁶. In the 2010s, the limit is set at 10 millions simultaneous connections²⁷. With the growing number of connected devices on the internet, concurrency is a very important property in the design of web applications. Moreover, the concurrency needs to be scalable to adapt to this growth of audience, as explained in the next paragraph.

Scalability The traffic of a popular web application such as Google search remains stable because of its popularity. The importance of the average traffic softens the

²⁴<http://blogs.cisco.com/news/cisco-connections-counter>

²⁵<http://www.internetlivestats.com/google-search-statistics/>

²⁶<http://www.kegel.com/c10k.html>

²⁷<http://c10m.robertgraham.com/p/manifesto.html>

occasional spikes. However, the traffic of a less popular web application is much more uncertain. For example, it might become viral when it is efficiently relayed in the media. The load of the web application increases with the growth of audience. The available resources needs to increase to meet this load. This growth can be steady enough to plan the increase of resources ahead of time, or it might be erratic and challenging. An application is scalable, if it is able to spread over resources proportionally as a reaction to the increasing growth of audience.

2.1.2.2 Time-slicing and Parallelism

Concurrency is achieved differently on hardware with a single or several processing units. On a single processing unit, the tasks are executed sequentially, interleaved in time. While on several processing units, the tasks are executed simultaneously, in parallel. Parallel executions uses more processing units to reduce computing time over sequential execution.

If the tasks are independent, they can be executed in parallel as well as sequentially. This parallelism is scalable, as the independent tasks can stretch the computation on the resources so as to meet the required performance.

However, the tasks within an application need to coordinate together to modify the application state. This coordination limits the parallelism and imposes to execute some tasks sequentially. It limits the scalability. The type of possible concurrency, sequential or parallel, is defined by the interdependencies of the tasks.

As explained in the previous section, Javascript requires a global memory to coordinate the execution of the callbacks. The event-loop is constrained within time-slicing concurrency to assure this coordination.

This thesis argues that there exists a parallel equivalent to the event-loop model, and that it is the pipeline execution model. The next section presents this parallel execution model.

2.1.2.3 Pipeline Execution Model

The pipeline software architecture is composed of isolated stages communicating by message passing to leverage the parallelism of a multi-core hardware architectures. It is well suited for streaming application, as the stream of data flows from stage to stage. Each stage has an independent memory to hold its own state. As the stages are independent, the state coordination between the stages are communicated along with the stream of data.

*



TODO
schema of a
pipeline

Each stage is organized in a similar fashion than the event-loop presented in section 2.1.1.3. It receives and queues messages from upstream stages, processes them one after the other, and outputs the result to downstream stages. The difference is that in the pipeline architecture, each task is executed on an isolated stage, whereas in the event-loop execution model, all tasks share the same queue, loop and memory store.

* The next section details further the incompatibility in their model and the resulting economical consequences. ⚠

TODO
schema to
compare the
event-loop
and the
pipeline,
and an
introducing
sentence.

2.2 An Economical Problem

With the rise of SaaS on the Web, the software industry are in charge of both the development and the execution of the software. The previous section presented these two aspects individually. This section presents the challenges encountered by conducting the two at such a large scale. It then focus on the subject and define the objectives of this thesis.

2.2.1 Disrupted Development

The economical context on the Web allows a project to grow from a very early stage to a large business. The economical constraints to meet are very different in the beginning and in the maturation of such project. In the early steps the constraints hold on the development. The project needs crucially to reduce development costs, and to release a first product as soon as possible. On the contrary, in the maturation of the project, the constraints hold on the performance. The product needs to be highly concurrent to meet the load of usage. The team needs to adapt to meet the different constraints, which implies a disruption in the evolution of the project. This section further details the reasons and consequences of this disruption.

2.2.1.1 Power-Wall Disruption

illustration:
heating
chipset
parallel
chipsets

Around 2004, the speed of sequential execution on a processing unit plateaued²⁸. Manufacturers reached what they called the power wall. They started to arrange transistors into several processing units to keep increasing overall performance while

²⁸<https://cartesianproduct.wordpress.com/2013/04/15/the-end-of-dennard-scaling/>

avoiding overheating problems. Therefore, the performance of the sequential execution plateaued as well. Parallelism is the only option to achieve high concurrency on this parallel hardware. But the isolation required by parallelism is in contradiction with the best practices of software development. This power-wall leads to a rupture between performance and maintainability.

2.2.1.2 Unavoidable Modularity

The best practices in software development advocate to gather features logically into distinct modules. This modularity allows a developer to understand and contribute to an application one module at a time, instead of understanding the whole application. It allows to develop and maintain a large code-base by a multitude of developers bringing small, independent contributions.

This modularity avoids a different problem than the isolation required by parallelism. The former intends to structure code to improve maintainability, while the latter improve performance through parallel execution. These two organizations are conflicting in the design of the application. The next paragraph presents the disruptions in the development of a web application implied by this conflict.

2.2.1.3 Technological Shift

Between the prototyping, and the maturation of a web application, the needs are radically different. During the initiation of a web application project, the economical constraint holds on the pace of development. The development reactivity is crucial to meet the market needs²⁹. The development team opt for a popular and accessible language to leverage the advantage of its community. It is only after a certain threshold of popularity that the economical constraint on performance requirements exceeds the one on development. The development team shifts to an organization providing parallelism.

This shift brings two risks. The development team needs to rewrite the code base to adapt it to a completely different paradigm. The application risks to fail because of this challenge. And after this shift the development pace slows down. The development team cannot react as quickly to user feedbacks to adapt the application to the market needs. The application risks to fall in obsolescence.

The risks implied by this rupture proves that there is economically a need for a solution that continuously follows the evolution of a web application. We present in the next section the proposition of this thesis for such a solution. It would allow

²⁹<https://www.cbinsights.com/blog/startup-failure-post-mortem/>

developers to iterate continuously on the implementation focusing simultaneously on performance, and on maintainability.

2.2.2 Seamless Web Development

This thesis is conducted in the frame of a larger work within the company Worldline : LiquidIT. Worldline identified that one of their need was to increase the time to market for its product. Worldline defines LiquidIT is defined as *a concept of flexible and cost-effective IT services that can be provisioned, built and configured in real time, allowing end-to-end financial transparency*. It precisely intends to provide *business agility, investment-free charging models, flexibility and ease of use*. The goal of this thesis in this larger work, is to allow the developer to focus solely on business logic, and leave the technical constraints of deployment to automated tools. This section presents the objective of this work to avoid the disruption in development, and provide a seamless development experience. Worldline develops and hosts real-time streaming Web services, as defined in the next paragraph.

2.2.2.1 Real-Time Streaming Web Services

This thesis focuses on web applications processing streams of requests from users in soft real-time. Such applications receive requests from clients through the HTTP protocol and must respond within a finite window of time. They are generally organized as sequences of tasks to modify the input stream of requests to produce the output stream of responses. The stream of requests flows through the tasks, and is not stored. On the other hand, the state of the application remains in memory to impact the future behaviors of the application. This state might be shared by several tasks within the application, and imply coordination between them.

The next section introduces the similarities and differences between the two programming models from the previous section. And then draws an equivalence. This equivalence is developed throughout this thesis.

2.2.2.2 Differences

Both paradigms encapsulate the execution in tasks assured to have an exclusive access to the memory. However, they provide two different models to provide this exclusivity resulting in two distinct programming models. Contrary to the pipeline architecture, the event-loop provides a common memory store allowing the best practice of software development to improve maintainability.

However, these two organizations are incompatible. Because of economical constraints, this incompatibility implies ruptures in the development. It represents additional development efforts and important costs. This thesis argues that it is possible to allow a continuous development between the two organizations, so as to lift these efforts and costs. This section presents the two programming models representing each an organization. Then it presents the possibility of an equivalence bridging the two. This equivalence is detailed further in the chapter 4 and 5 of this thesis.

2.2.2.3 Equivalence

With this equivalence, it would be possible to express an application following the design principles of software development, hence maintainable. And yet, the execution engine could adapt itself to any parallelism of the computing machine, from a single core, to a distributed cluster. Because of the equivalence between these two models, the development team could iterate testing the two models for their different concerns about the implementation : performance and maintainability. In the beginning of a project, the team focuses on maintainability and evolution, discarding the scalable performance concerns. And as the project gather audience and the performance concerns become more and more critical, the development team progressively take into account this performance concerns.

The goal of conciliating these two concerns is not new. The next chapter presents all the results from previous works needed to understand this work, up to the latest results in the field.

This thesis proposes to provide an equivalence between the two memory models for streaming web applications. The next section further describes the similarities and differences between the two models. The equivalence would allow a compiler to transform an application expressed in one model into the other. With such a tool, a development team could rely on the common memory store of the event-loop execution model, and focuses on the maintainability of the implementation. And continuously compiles during the development the event-loop implementation to the pipeline architecture to assure that the execution can be distributed on a parallel architecture.

Chapter 3

Software Design, State Of The Art

Contents

3.1 Software Maintainability	24
3.1.1 Modular Programming	25
3.1.1.1 Design Choices	25
3.1.1.2 Programming Models	26
3.1.1.3 Performance Limitations	28
3.1.2 Performance Improvements	29
3.1.2.1 Concurrent Programming	30
3.1.2.2 Compilation	32
3.1.2.3 Accessibility Limitations	35
3.2 Software Performance	35
3.2.1 Parallel Programming	36
3.2.1.1 Asynchronous and Isolated Process Parallelism .	37
3.2.1.2 Stream Processing Systems	38
3.2.1.3 Elitism of Parallel Programming	39
3.2.2 Maintainability Improvements	40
3.2.2.1 Execution Organization	41
3.2.2.2 Memory Abstraction	42
3.2.2.3 Lack of Higher-Order Programming	42
3.3 Seamless Development	45

3.3.1	Equivalence	47
3.3.1.1	Rupture Point	47
3.3.1.2	Invariance	48

“A designer is responsible for producing the greatest benefit for any given investment of time, talent, money, and other resources.”

—K. Sullivan, W. Griswold, Y. Cai, B. Hallen [**Sullivan2001a**]

* The growth of the web and Software as a Service (SaaS) revealed the importance of previously unknown economic constraints. The same company carries both development and exploitation of a service in scale of unprecedented size. Development costs are reduced by following best practice, and by building maintainable softwares. However, as seen in the previous chapter, it is compensated by increasing hardware performance, which eventually rises exploitation costs. Similarly, exploitation costs are reduced by following more efficient programming models. But again, it rises development costs. So a SaaS company needs to cleverly allocate its budget between development and exploitation so as to limit the overall cost.

⚠
TODO
the next
paragraph
expresses
well the idea,
but is not
clear enough
yet

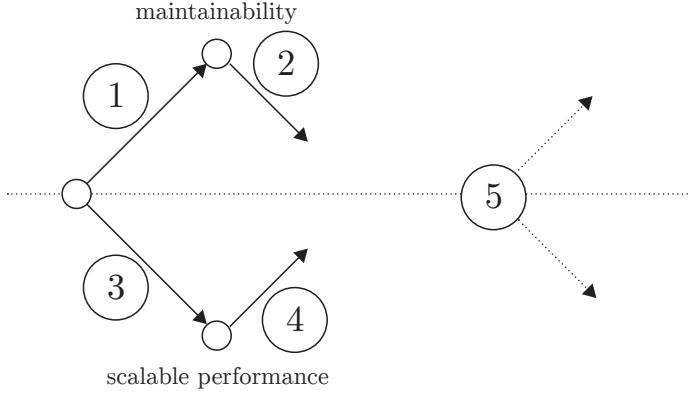
Eventually, a company faces the problem of scalability limitations. Compensating development with hardware becomes unsustainable. The company has no choice but to commit huge development efforts to get correct performances. This chapter draws a broad view of the relation between the orientation of development toward maintainability or scalable performance, and its consequences.

The best practices in software design advocate to decompose a problem into many subproblems. The decomposition of the implementation of a problem improves directly its maintainability, development scalability and evolution. This chapter present some of these best practices, *e.g.* modular programming, structured design [73], hierarchical structure [23] and object-oriented programming.

A few decades ago, the best practices were not concerned with execution performance. Moore’s law [57] was wrongly interpreted as the assurance that hardware could always increase execution speed. But eventually, the clock speed of processors plateaued ¹[Bohr2007], and the processing units were organized as several execution units to continue improving performances. But this hardware improvement could not anymore increase the execution speed without any additional development effort.

The best practices of software design then inherited two goals : to assure a scalable implementation evolution by decomposing it into subproblems, as well as assure a scalable parallel execution by decomposing the execution onto the several execution units. As D. L. Parnas showed in 1972 [63], it seems challenging to develop a software following a decomposition that satisfies both goals.

¹<https://cartesianproduct.wordpress.com/2013/04/15/the-end-of-dennard-scaling/>



The evolution of economic constraints often force an SaaS company to switch from scalable development to scalable performance. This switch implies huge development efforts. There has been many attempts at reconciling the two goals to reduce these costs. But none seems really convincing enough to be widely adopted.

The schema above is a graphical representation of the organization of this chapter, and more generally of the state of the art of software design. The focus on development scalability, noted by number 1, are addressed in section 3.1 while the focus on performance scalability, noted by number 3, are addressed in section 3.2. Each of these direction of development contains works trying to meet the requirements from the opposing category, noted by number 2 and 4. And finally, section 3.3 presents the objectives for this thesis, noted by number 5.

3.1 Software Maintainability

“It is becoming increasingly important to the data-processing industry to be able to produce more programming systems and produce them with fewer errors, at a faster rate, and in a way that modifications can be accomplished easily and quickly.”

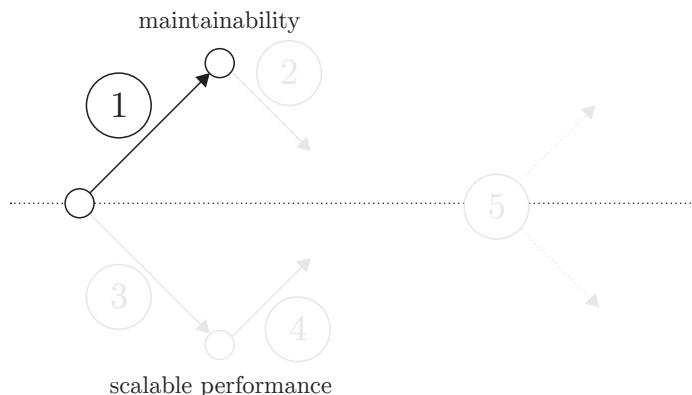
—W. P. Stevens, G. J. Myers, L. L. Constantine [73]

In order to improve and maintain a software system, it is important to hold in mind the mental representation behind its implementation. Architects, and mechanical engineers draw codified plans to share their mental representations with peers and building teams. software design is an exception in that the implementation is both the shared mental representation, and the actual product. The mental representation is often lost in technical details and optimizations for the actual product.

This problem becomes even more critical as the system grows in size. Therefore, it is crucial to decompose the system into smaller subsystems easier to grasp individually. This section shows the theories, programming languages and frameworks helping this decomposition.

3.1.1 Modular Programming

The modularity of a software implementation is about enclosing the subproblems and bringing the relevant interfaces to allow these parts to be composed. It allows greater design to emerge from the composition of smaller components. Such modularity helps organizing the implementation to reflect the underlying mental organization. The modularity in software design improves the maintainability of an implementation, as presented in the following schema. It allows to limit the understanding required to contribute to a module [73]. And it reduces development time by allowing several developers to simultaneously implement different modules [Cataldo2006, 79].



The section 3.1.1.1 is about the decomposition of a problem into subproblems. Then, the section ?? is about the bringing the interfaces allowing their composition. Finally, the section ?? presents the consequences of this decomposition on performance.

3.1.1.1 Design Choices

* In the decomposition of a large problem into smaller subproblems, there are two design choices. The first one is the granularity, and organization of the subproblems within the system decomposition. The second one is the organization of the implementation within the subproblems to improve maintainability.



TODO this introduction is not very clear

illustration:
spaghetti
programming

System Decomposition Dijkstra firstly developed the concept of Structured Programming [Dijkstra1970], which later led to modular programming. Structured programming is about drawing clear interfaces around a piece of implementation so that the execution is clearly enclosed inside. At a fine level, it helps avoid spaghetti code [22], and at a coarser level, it structures the implementation [23] into modules, or layers. The next paragraph explains further the criteria to draw the borders around modules.

illustration:
lasagna pro-
gramming

Decomposition Criteria * The criteria to decompose the system into well defined modules are coupling and cohesion [73]. The coupling defines the strength of the interdependence between modules, while cohesion defines how strongly the features inside a module are related. Low coupling between modules and high cohesion inside modules helps logically organize, and understand the implementation. Hence, it improves its maintainability. The next paragraph presents the approach to build modules helping with the evolution of the implementation.

⚠
move this
paragraph
closer to
OOP ?

Development Evolution To improve maintainability of implementation, the modular organization should isolate the evolution of a module from impacting the rest of the implementation. The Information Hiding Principle [63], and the Separation of Concerns [Tarr1999, Hursch1995] are two similar approaches to do so. The information hiding principle advocates to encapsulate a specific design choice in each module. The Separation of Concerns advocates each module to be responsible for one and only one specific concern. Examples of separation of concerns are the separation of the form and the content in HTML / CSS, or the OSI model for the network stack.

3.1.1.2 Programming Models

The previous section presented the design choices to build modules. This section presents the programming models providing the interfaces to glue the modules together. It focuses on two main programming models currently used in the industry, object oriented programming and functional programming.

illustration:
multiple cells
communicat-
ing

Object Oriented Programming Alan Kay, who coined the term, states that Object Oriented Programming (OOP) is about message-passing, encapsulation and late binding. (There is no academic reference for that, only a public mail exchange².)

²http://userpage.fu-berlin.de/~ram/pub/pub_jf47ht81Ht/doc_kay_oop_en

Message-passing and late binding loosen coupling between objects, while encapsulation is intended to increase cohesion ³.

The very first OOP language was Smalltalk [Goldberg1984]. It defined the core concept of OOP. Nowadays, the major emblematic figures of OOP in the software industry are C++ and Java [Gosling2000, Stroustrup1986]. Though, the trend seems to digress from these languages to evolve toward a more dynamic approach, closer to Functional Programming. Indeed Javascript adopts some functional features such as dynamic typing and higher-order functions [Ecma1999].

Functional Programming The definition of pure Functional Programming resides in manipulating only mathematical expressions - functions - and forbidding state mutability. The absence of state mutability makes a function referentially transparent, and thus side-effect free. The most important pure Functional Programming languages are Scheme [Rees1986], Miranda [Turner1986], Haskell [Hudak1992], Erlang [JoeArmstrong] and Standard ML [Milner1997].

* However, the functional programming concepts are also implemented in other languages along with mutable states. Major imperative programming languages now commonly present higher-order functions and lazy evaluation to help loosen the couple between modules, define more generic and reusable modules. *In fine*, it helps developers to write applications that are more maintainable, and favorable to evolution [Hughes1989, Turner1981].



TODO find a better argument to say that immutability is unadapted to every day programming

Higher-Order Programming *



Higher-order programming allows to manipulate functions like any other primary value : to store them in variables, or to pass them as arguments. It replaces the need for most modern object oriented programming design patterns ⁴. For example Inversion of Control [Johnson], the Hollywood Principle [Sweet1985], and Monads [Wadler1992]. Higher-order programming help loosen coupling, thus improve maintainability.

If possible, include this reference : Continuations and coroutines [35]

In languages allowing mutable state, higher-order functions are implemented as closure, to preserve the lexical scope [74]. A closure is the association of a function and a reference to the lexical context from its creation. It allows this function to access variable from this context, even when invoked outside the scope of this context. It eventually tangles the memory references so that it requires a global memory.

³<http://williamdurand.fr/2013/06/03/object-calisthenics/>

⁴<http://stackoverflow.com/a/5797892/933670>

Lazy Evaluation Lazy evaluation is an evaluation strategy allowing to defer the execution of a function only when its result is needed. The lazy evaluation of lists is equivalent to a stream with a null-sized buffer, while the opposite, eager evaluation, corresponds to an infinite buffer [VanRoy2003]. Indeed, the dataflow programming paradigm resulting from lazy lists is particularly adapted for stream processing applications.

The lazy evaluation, as well as streams are powerful tools for structuring modular programs [Sussman1983].^s Lazy evaluation allows the execution to be organized as a concurrent pipeline, as the stages are executed independently for each element of the stream. But this concurrency requires immutability of state, or at least isolation of side-effects. The next section addresses the consequences of higher-order programming and lazy evaluation on parallelism.

3.1.1.3 Performance Limitations

??

Functional programming greatly support modularity to improve the maintainability of an application, and its resilience to evolution. However, the closures introduced by higher-order programming require to share the execution context among modules. The previous chapter show that sharing makes parallelism difficult. It is the reason why that maintainability and performance seem hardly compatible. This section explore in further details the limitation of modular programming regarding performances.

Tighten Memory Closures are implemented in languages using a global memory. And by exchanging closures, two modules intricately share their contexts of execution. Higher-order programming loosen the couple on the implementation level of the modules, but tighten it on the execution level. It improves modularity, but it inherently worsens the parallelization hence the performance scalability.

Scalability Limitations Parallelizing the execution increases the performances [7, 32] But the parallelism is limited because the execution portions sharing state need to be scheduled sequentially. Hence, to increase the parallelism and performance, the concurrent executions need to be independent, or to coordinate to be scheduled sequentially [34, 33, 60, 31]. We explain further the reasons of these limitations, and the improvement solutions in the next section.

* The modular organization of implementation is opposed to the organization ⚠

TODO the
transition is
not very clear

	Model	Examples	Higher-Order Programming	Modularity Closures	Lazy Evaluation / Streams
Modular Programming	Object-Oriented Programming	C/C++, Java	V	X	X
	Functional Programming	Javascript	V	V	V

Table 3.1 – Synthesis of the state of the art in modular programming

favoring the parallelization of the execution. The former organization supports the development scalability, while the latter supports performance scalability. A program cannot trivially follow an organization that support both development evolution, and performance. However, D. Parnas advocates the use of an assembler to conciliate the two approaches [63]. The next section shows the improvements for performance and parallelism . Then, section 3.2 shows the techniques for parallelism.

3.1.2 Performance Improvements

To assure its integrity, the global state of an application imposes the concurrent executions to coordinate their accesses. This coordination is responsible of the atomicity, and exclusivity of the accesses. It assures the invariance of the state during its atomic manipulation. So that developers can group operations in atomic manipulations so as to avoid corruption of the state.

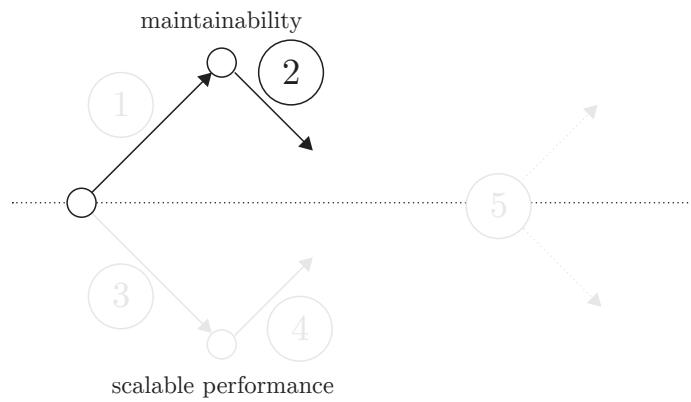
The invariance is assured differently depending on how the state is shared among the concurrent execution. To increase performance, concurrent executions needs to be as independent as possible to be executed in parallel ⁵.

Isolation If different concurrent executions are commutative [68, 17], and they share no portion of the state and can be isolated and executed in parallel.

Immutability Otherwise, the sharing state portions needs to be immutable to conserve invariance and parallel execution [Gordon2012, 54].

⁵<http://joeduffyblog.com/2010/07/11/thoughts-on-immutability-and-concurrency/>

Synchronization If different concurrent executions needs mutation on the state, their accesses are scheduled sequentially.



The next few paragraphs presents the different models to assure invariance in concurrent execution, while conserving modular programming, as illustrated in the schema above. Section 3.1.2.1 presents the programming models providing synchronization and immutability for concurrent executions. Section 3.1.2.2 presents compilation methods to parallelize sequential programs.

3.1.2.1 Concurrent Programming

Concurrent programming provides to the developer the mechanisms for concurrent execution, while conserving a global memory model.

illustration:
feu rouge et
rond point

There are two scheduling strategies to execute tasks sequentially on a single processing unit, cooperative scheduling and preemptive scheduling. Cooperative scheduling allow a concurrent execution to run until it yields back to the scheduler. Each concurrent execution has an atomic, exclusive access on the memory. On the other hand, preemptive scheduling allows a limited time of execution for each concurrent execution, before preempting it. It assures fairness between the tasks, such as in a multi-tasking operating system, but as the preemption happens unexpectedly, the developer needs to lock the shared state to assure atomicity and exclusivity. The next paragraphs presents the event-based programing model, based on cooperative scheduling, and the multi-threading programming model, based on preemptive scheduling.

Event-Driven Programming Event-driven programming explicitly queues the concurrent executions needing access to shared resources. The concurrent executions are scheduled sequentially to assure exclusivity, and cooperatively to assure atomicity.

As presented in the previous chapter, web servers need to be highly concurrent, and efficient. The event-driven model is very efficient to serve websites, as it avoids contention due to waiting for shared resources like disks, or network. Web servers like Flash [62], Ninja [30] thttpd⁶ and Nginx⁷ were designed following this model. However, a drawback of this model was that the execution context is lost at each event. The developer needs to explicitly transfer the relevant state to continue the execution from one event execution to another.

Cooperative threads, or fibers, addressed this drawback [4, 12]. The execution is not ripped into several events. It yields and resume exactly at the same point after the completion of an asynchronous operation, conserving its context. However, the developer needs to be well aware of the asynchronous calls to assure the atomicity⁸.

The problem of losing the execution context disappears with closures in higher-order programming. Moreover, the continuation passing style used in higher-order programming requires the developer to be aware of the asynchronous rupture in the execution, so as to assure atomicity [74]. And because an asynchronous call doesn't wait for the completion of the operation, the asynchronous control flow is not limited to be linear like in threads. Multiple asynchronous calls are made in parallel. Several execution models proposed this event-based programming model, like TAME [48], Node.js⁹ and Vert.X¹⁰.

However, as the shared memory is global and all the execution portions need atomic access, they are not parallel, but sequentially concurrent. The next paragraph presents the multi-threading and associated synchronization mechanisms to try to improve the parallelism of execution using finer granularity of atomic execution and exclusivity.

Multi-Threading Programming Threads are light processes sharing the same memory execution context within an isolated process [23]. They wait for completion of each operation, and are preemptively scheduled to avoid blocking the global progression. This preemption breaks the atomicity of the execution, and the parallel execution breaks the exclusivity. To restore atomicity and exclusivity, hence

⁶<http://acme.com/software/thttpd/>

⁷<https://www.nginx.com/>

⁸<https://glyph.twistedmatrix.com/2014/02/unyielding.html>

⁹<https://nodejs.org/en/>

¹⁰<http://vertx.io/>

assure the invariance, multi-threading programming model provide synchronization mechanisms, such as semaphores [**Dijkstra**], guarded commands [21], guarded region [**Hansen1978a**] or monitors [40]. They assure an execution region to have exclusive access over a cell of the global state.

“The purpose of explicit synchronization is to manage the timing of side-effects in the presence of parallelism.”

—Chris Quenelle¹¹

Developers tend to use the global memory extensively, and threads require to protect each and every shared memory cell. This heavy need for synchronization leads to bad performances, and is difficult to develop with [4]. The next paragraph present work intending to improve performance by reducing the lock granularity to a minimum.

Lock-Free Data-Structures The wait-free and lock-free data-structures reduce the exclusive execution to a few atomic operations [**Lamport1977**, **Herlihy1988**, **Herlihy1990**, **Herlihy1991**, **Anderson1990**]. They are based on transactional memories [**Harris2010**], which provide atomic read and write operations on a shared memory. Lock-free data-structures arrange these atomic operations so as to keep invariance without the need to lock. They provide concurrent implementation of basic data-structures such as linked list [**Valois1995**, **Timnat2012**], queue [**Sundell2003**, **Wimmer2015**], tree [**Ramachandran2015**] or stack [**Hendler2004**].

However, even if they are theoretically infinitely scalable, they are hard to come with, and are not fit for every problem.

This section showed that it is difficult for developers to assure the invariance of memory in the context of parallel programming. Multi-threading programming is inefficient and difficult to program with. Event-driven programming is easy to develop with but limits parallel execution to assure exclusivity. The global memory requires the synchronization between concurrent executions to some extent.

Synchronization inherently limits the scalability. The next section presents compilation methods to improve parallelism, by extracting immutability and isolation from sequential programs.

3.1.2.2 Compilation

“It is a mistake to attempt high concurrency without help from the compiler”

—R. Behren, J. Condit, E. Brewer [11]

When showing the incompatibility between the two organization, D. Parnas advocated conciliating the two methods using an assembler to transform the development organization into the execution organization [63]. This section presents the state of the art to extract parallelization from sequential programs through code transformation and compilation.

Parallelism Extraction As the only requirement to parallelism is the commutativity of operations [68, 17], a compiler needs to identify the commutative operations transform a sequential program so as to parallelize its execution [68].

An important work was done to parallelize loop iterations [Mauras1989, 6, 10, 66], particularly using the polyhedral compilation method [Yuki2013, Grosser2011, Trifunovic2010, Bastoul2004]. However, this data parallelism is limited to scientific applications because of their heavy use of loops on matrices and vectors. The performance gains are limited in common sequential programs, as the execution remains sequential outside of loops [7, 17].

To improve performance gains further, some compilers identify the data-flow inside sequential programs to allow pipeline parallelism on the whole program, and not only on its loops [Beck1991, Li2012]. Moreover, the data-flow representation and execution of a program is well suited for modern data processing applications [24], as well as web services [69]. *

However, the limitation of modular programming regarding parallelization persists. In a purely functional language with immutability, higher-order functions are referentially transparent which implies commutativity hence parallelism *. However, in a functional language with mutable data, closures remains a challenge to parallelize, because of the memory references shared across the program [Harrison1989, Nicolay2010, 54]. The next two paragraphs presents two directions to improve the state of the art in parallel compilation. The first paragraph presents static analysis, while the second presents annotations systems.

⚠
TODO
Extract
pipeline
parallelism
compilers
from : Load
balanced
pipeline
parallelism
[Kamruzzaman2013]

⚠
Add reference
of parallel
purely
functional
languages

Static analysis Compilers analyze the control-flow of a program to detect the side-effects causing dependencies between statements [Allen1970]. The point-to analysis, presented by L. Andersen [8] is a popular approach to identify these side-effects in the memory representation. The points-to-analysis was adapted for Javascript [43, 72, 77], and is a useful tool to analyze a program. However, this analysis is not sufficient to track the dynamic control-flow of higher-order functions [Shivers1991] like used in Javascript.

The Operational Semantics is an example of abstract interpretation technique that allows to statically reason on the behavior of programs[51, 71, 28, 27, 13].

* Abstract interpretation techniques are more adapted for program with higher-order functions, and are successfully used for security applications [**Chudnov2015**, **Dolby2015**, 41, 44, 80, 52]*.



TODO
review this
paragraph

However, static analysis techniques are too imprecise, and expensive for the performance gain to be profitable in languages as dynamic as Javascript. Instead, some compilers relies on annotations from the developers.



Update the
citation for
Dolby2015

Annotations Extracting parallel dataflow from an imperative, sequential implementation is a hard problem [**Johnston2004a**]. Some works proposed to rely on annotations from the developer to help the identify the possible side-effects between operations [**Vandierendonck2010a**, 24].

Many compilers rely on annotations from the developer to build highly parallel executables. Such annotations are especially relevant for accelerators such as GPUs or FPGAs, because the development effort yield huge performance improvements. Examples of such compilers are OpenMP [**Dagum1998**], OpenCL [**Stone2010**], CUDA [**Nvidia2007**] Cg [53], Brook [14], Liquid Metal [**Huang2008**].

However, the burden of detecting commutativity of operations, or independence of operations fall back to the developer. In this regard, these solutions successfully improve performances, but are unable to fix the rupture between performance and maintainability. These solutions are indeed very close to the performance oriented solutions presented in the section 3.2.

Compilation Limitations The static analysis of static, low level languages like FORTRAN or C, brings performance improvements. However for more dynamic, higher-level languages like Javascript, the static analysis is not sufficient to identify correctly the dependencies between operations to parallelize them. And parallel compilers often fall back on relying on annotation provided by developers. So, in this regards, it seems that the accessibility of development gained by higher-level programming is detrimental to performance.

3.1.2.3 Accessibility Limitations

The two previous sections showed that parallel programming is difficult, and that compilation is not yet mature enough to lift this burden from the developer.

Indeed, preemptive scheduling and synchronization mechanisms are known to be hard to manage by developers, and to impact performances negatively. On the other

	Model	Examples	Concurrency	Parallelism	
			Synchronization	Immutability	Isolation
Concurrent Programming	Event-driven	Node.js, Vert.X, TAME	V	x	x
	Multi-Threading	Lock, Mutex	V	x	x
	Lock-Free		V	x	x
	Data Structure		V	x	x
	Compilation	Loop parallelization CUDA, OpenCL, OpenMP ...	V	V	V

Table 3.2 – Synthesis of the state of the art in performance improvement of modular programming

hand, cooperative scheduling provides a more accessible invariance abstraction for concurrent programming, but seems limited to sequential execution because of it.

The next section presents the programming paradigms focusing on performance more than development accessibility, and then presents the works to improve accessibility.

3.2 Software Performance

Moore's law [57] which forecasts the density of transistors per processing unit, was wrongly interpreted to promise the exponential evolution in the sequential performance of the processing unit, and the assurance for the software industry of always faster hardware. But as transistors attained a critical size, the reduction in power required by transistor predicted by the Dennard's MOSFET scaling [**Dennard2007**] stopped¹². The ever growing number of transistor predicted by Moore's law are arranged in parallel architecture to continue increasing the performance of processing units. Parallel programming became the only solution for scalable performance, at the expense of development effort.

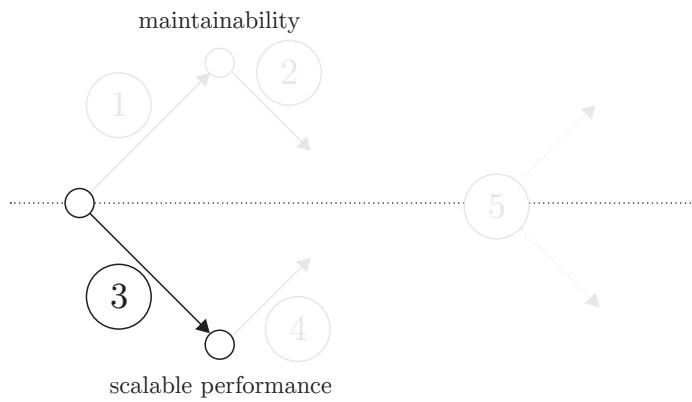
This section presents the parallel programming solutions and their limitations in

¹²<https://cartesianproduct.wordpress.com/2013/04/15/the-end-of-dennard-scaling/>

accessibility, and then the improvements to overcome these limitations.

3.2.1 Parallel Programming

Concurrent programming is based on the causal ordering of execution. The ordering of operations is local within a synchronous execution, while the concurrent executions are causally ordered. It leads to parallel execution with some coordinations such as synchronization, immutability or isolation. As Lamport showed [49], and Reed related later [67], This causal order is sufficient to execute correctly a system in parallel, such as in distributed system.



This section presents first hand the theoretical and programming model based on asynchronous communication and isolated execution for parallel programming, as illustrated on the schema above. It then continues with stream processing programming model. And finally, it concludes on the limitations of parallel programming regarding accessibility.

3.2.1.1 Asynchronous and Isolated Process Parallelism

The Flynn's taxonomy [Flynn1972] is the most commonly used to categorize parallel execution. It separates the flow of instructions, and the flow of data ; each being unique, or multiple. All the current parallel programming model currently belong to the category Multiple Instruction Multiple Data (MIMD), which is further divided into Single Program Multiple Data (SPMD) [Auguin1983, Darema1988, Darema2001] and Multiple Program Multiple Data (MPMD) [Chang1997, Chan2004]. MIMD implies several threads of execution processing several stream of data.

*

The difference between SPMD and MPMD is in the representation of the execution in implementation. SPMD organizes the implementation as a single execution replicated on many processing units. While MPMD organizes explicitly the different threads of execution in the implementation. Examples of SPMD programming languages are Split-C [Culler], CRL [Johnson1995] and Composite C++ [K.ManiChandy2005]. Examples of MPMD programming languages are Mentat [Grimshaw1991], Fortran M [Foster1995b] and Nexus [Foster1996]. SPMD is close to the model presenting parallel improvements over modular programming presented in section 3.1.1.2. While MPMD is closer to the programming models based on isolated process presented in the remaining of this section. The coordinations between these threads of execution were done by message passing, using PVM [Sunderam1994], MPI [Snir1996, Walker1996], SOAP, or the more recent REST protocols.

Theoretical Models The communication in reality are subject to various faults and attacks [Lamport1982] and too slow compared to execution to be synchronous. The Actor model is one of the first programming model to be explicitly designed to take these physical limitations in account [Hewitt1977a]. It allows to express the computation as a set of communicating actors [Clinger1981, 38, 37]. In reaction to a received message, an actor can create other actors, send messages, and choose how to respond to the next message. All actors are executed concurrently, and communicate asynchronously. An asynchronous communication implies that the sender continues its execution immediately after sending the message, before receiving the result of the initiated communication.

In the Actor Model, everything is an actor, even the simplest types like numbers. This level of granularity is unachievable in practice due to overhead from the asynchronous communications. Most implementations adopt a granularity on the process or function level.

Coroutines are autonomous programs which communicate with adjacent modules as if they were input and output subroutines [19]. It is the first definition of a pipeline to implement multi-pass algorithms. Similar works include the Communicating Sequential Processes (CSP) [Brookes1984, 39], and the Kahn Networks [45, 46].

Programming Languages The theoretical models presented above are implemented in industrial languages such as Akka Scala and Erlang.

Scala is an attempt at unifying the object model, and functional programming

[Odersky2004]. Akka¹³ is a framework based on Scala, following the Actor model to build highly scalable and resilient applications. Play¹⁴ is a web framework based on top of Akka.

Erlang borrows the Actor model as well. It is a functional concurrent language designed by Ericsson to operate telecommunication devices [JoeArmstrong, Nelson2004]

* The field of concurrent programming is so vast it is impossible to relate here every programming languages. The previous examples are only the best known. The next focus focuses on streaming real-time applications.



review this paragraph and the transition to the next section

3.2.1.2 Stream Processing Systems

All the solutions previously presented are designed to build general distributed systems. In the context of the web, a real-time application must process high volumes streams of requests within a certain time. Because these systems are key to business, their reliability and latency are of critical importance. Otherwise, input data may be lost or output data may lose their value. These requirements are challenging to meet in the design of such system.

Data-stream Management Systems Database Management Systems (DBMS) historically processed large volume of data, and they naturally evolved into Data-stream Management System (DSMS) to process data streams as well. Because of this evolution, they are in rupture with imperative languages presented until now, and borrow the syntax from SQL.

DSMS concurrently run SQL-like requests on continuous data streams. The computation of these requests spread over a distributed architecture. Among the early works, we can cite NiagaraCQ [16, 59], Aurora [1, 3, 9] which evolved into Borealis [2], AQuery [Lerner2003], STREAM [Arasu2003, Arasu2005] and TelegraphCQ [47, 15]. More recently, we can cite DryadLINQ [42, 81], Apache Hive [Thusoo2009]¹⁵, Timestream [65] and Shark [Xin2013].

Pipeline Architecture As presented in the previous section, streaming and lazy-evaluation composition both allow a loosely coupled yet efficient composition. The pipeline architecture takes advantage of this, and composes the parallel execution in a stream, the output of one feeding the input of the next.

¹³<http://akka.io/>

¹⁴<https://www.playframework.com/>

¹⁵<https://hive.apache.org/>

SEDA is a precursor in the design of pipeline-based architecture for real-time web applications [78]. It organizes an application as a network of event-driven stages connected by explicit queues. The event-driven paradigm is similar to previous web servers implementations like Ninja and Flash [30, 62]. SEDA improves with the pipeline organization in stages.

Several projects followed and adapted the principles in this work. StreaMIT is a language to help the programming of large streaming application [75]. Storm [76] is designed by and used at Twitter to process the heavy streams of tweets. Among other works, in the industry, there are CBP [50] and S4 [61], that were designed at Yahoo, Millwheel [5] designed at Google and Naiad [Murray2013] designed at Microsoft.

In the litterature, there are Spidle [18], Pig Latin [Olston2008], Piccolo [64], Comet [36], Nectar [Gunda2010], SEEP [56] and SDG [24]

Transition on the limitations of software parallelism

3.2.1.3 Elitism of Parallel Programming

In these parallel programming models describing a network of isolated execution, the topology of the network is statically defined. The dynamical modification of the topology is impossible. It is not possible to dynamically manipulate execution containers, like it is possible to manipulate functions. Therefore, higher-level programming is impossible, so modular programming is limited.

Moreover, the decomposition of the execution and memory required is difficult for most developers to manage correctly. * It implies to keep two mental representation of the implementation. One for the modularity, and one for the decomposition of execution. Parallel programming remains hard, and is accessible only to an elite of developers.

The next section presents different improvements to make parallel programming more accessible to common developers.



TODO make sure this idea is correctly developed throughout this chapter : parallel programming implies to keep a double mental representation.

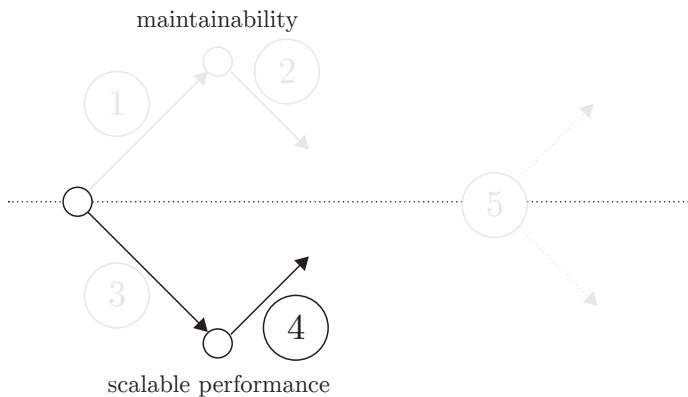
3.2.2 Maintainability Improvements

Introduction

As the limitations above state, the difficulties in parallel programming are the absence of higher-level programming, and the decomposition of execution and memory required for parallelism. This section presents the improvements on accessibility in parallel programing, as illustrated in the schema below.

	Model	Examples	Concurrency Synchronization	Parallelism	
				Immutability	Isolation
Parallel Programming	Actor Model	Scala, Akka,	x	V	V
		Play	x	V	V
		Erlang	x	V	V
	Stream Processing (Timestream)	DSMS	x		
		Pipeline	x	V	V
		(SEDA)			

Table 3.3 – Synthesis of the state of the art in parallel programming



These improvements holds on execution or memory. This section presents firstly the improvements to ease the decomposition of execution through design patterns and different granularities. It then presents the improvements to helps development in face of the required memory isolation for distribution. It then transitions to the lack of improvements for higher-order programming in parallel programming.

3.2.2.1 Execution Organization

All of the parallel programming models presented above decompose the execution into isolated parallel executions, like actors. This decomposition is difficult as it implies to manage both the parallel decomposition of execution and the modular decomposition of implementation. Most developers are unable to manage efficiently the two decompositions. The next paragraphs presents some solutions to mitigate this duality. *



TODO weak argumentation

Design Patterns To reduce the difficulties of the decomposition of the execution into actors, algorithmic skeletons propose predefined patterns that fit certain type of problems [Cole1988, Gonzalez-Velez2010, 20, 55]. A developer implements the problem as a specific case of a skeleton. It simplifies the communications, so that the developer can focus on its problem independently of message passing required by the distribution of execution.

Granularity The Service Oriented Architectures (SOA), and more recently Microservice[Namiot2014, Fowler2014, Namiot2014, 25] allow developers to express an application as an assembly of services connected to each others. Some examples of frameworks are OSGi¹⁶, EJB¹⁷, Spring¹⁸, and Seneca¹⁹. It intends to adjust the granularity of execution decomposition to help developers to fit the two organizations, the modular organization and the parallel execution organization [Adam2008].

Dynamic Distribution An interesting work following SEDA, is Leda [69, 70]. It follows the PCAM design methodology [Foster1995] to propose a model where the stages of the pipeline are defined only by their role in the application. The actual execution distribution in stages is defined automatically, only after the development, during deployment. This automation blurs the distinction between the parallel organization of execution, and the modular organization of implementation. It manages the execution organizations to helps the developer focus on the modular organization.

These works helps developers to decompose and distribute the execution of a parallel application. However, they still propose a distributed memory model. It doesn't allow higher-level programming, and it requires developers to distribute the state of the application, and assure its isolation. For these two reason it is difficult to manage. The next paragraph presents works to improve on the distribution of memory.

3.2.2.2 Memory Abstraction

Parallelization of the execution eventually requires the distribution of the memory. The Partitioned Global Address Space (PGAS) provides the developers with a uniform memory access on a distributed architecture. It attempts to combine the advantage of SPMD programming style for distributed memory systems, with the

¹⁶<https://www.osgi.org/developer/specifications/>

¹⁷<http://www.oracle.com/technetwork/java/javaee/ejb/index.html>

¹⁸<http://projects.spring.io/spring-framework/>

¹⁹<http://senecajs.org/>

	Model	Examples	Higher-Order Programming	Modularity Closures	Lazy Evaluation / Streams
Maintainability Improvements	Algorithmic Skeletons SOA, Microservices Dynamic Distribution PGAS	Map Reduce EJB, Spring, ... Seneca LEDA Chapel, X10 ...	x x x x	x x x x	V V V x

Table 3.4 – Synthesis of the state of the art in maintainability improvements for parallel programming

data referencing semantics of shared memory systems. Each computing node executes the same program, and provide its local memory to be shared with all the other nodes. The PGAS programming model assure the remote accesses and synchronization of memory across nodes, and enforces locality of reference, to reduce the communication overhead. Examples of implementation of the PGAS model are Chapel [Chamberlain2007], X10 [Charles2005]. Unified Parallel C [El-Ghazawi2006], CoArray Fortran [Numrich1998] and OpenSHMEM [Chapman2010].

The PGAS programming model is similar to the Multi-Threading Programming model presented in section 3.1.2.1. However, it is more focused on performance than modularity. It focuses on scientific applications with intensive computing such as matrices multiplication.

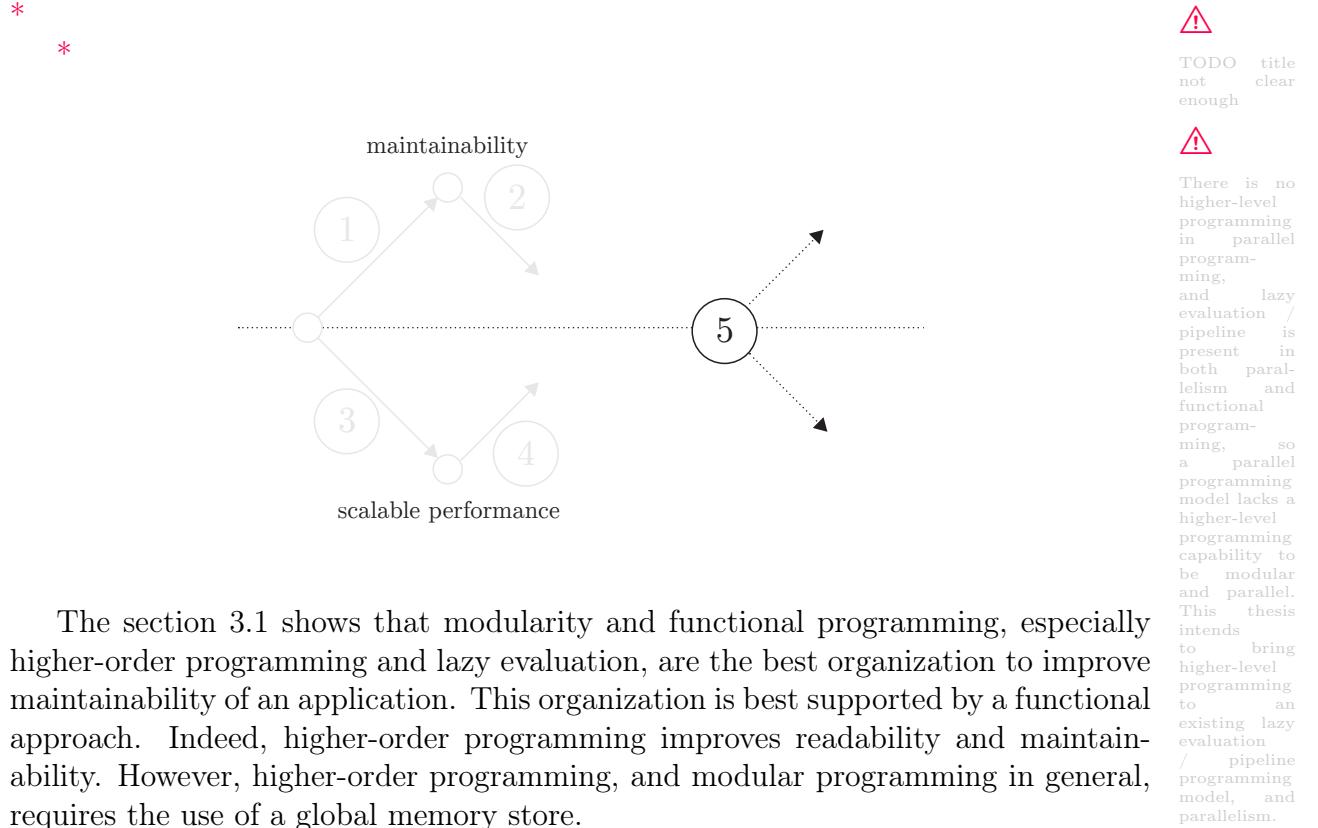
3.2.2.3 Lack of Higher-Order Programming

The programming models of this section lack higher order programming. As showed earlier in section 3.1.1.2, higher-order programming is important for modular design and maintainability of the implementation. In this regards, parallel programming seems currently incompatible with modular programming. The next section presents the proposition of this thesis to bring parallel programming to a higher-order programming language.

				Modularity		Concurrency		Parallelism	
				Closures	Lazy Streams	Synchronization	Immutability	Isolation	
Model	Examples	Higher-Order Programming	Evaluation / Streams						
Modular Programming	Object-Oriented Programming	C/C++, Java	V	X	X	X	X	X	X
	Functional Programming	Javascript	V	V	V	X	X	X	X
	Event-driven		Node.js, Vert.X, TAME Lock, Mutex	V	V	?	V	X	X
	Multi-Threading		?	?	V	?	X	X	X
	Lock-Free Data		?	?	V	?	X	X	X
	Structure		Loop parallelization	X	X	V	X	X	X
	Compilation		CUDA, OpenCL, OpenMP ...	X	X	V	V	V	V
	Concurrent Programming		Actor Model	Scala, Akka, Play, Erlang, DSMS	X	V	X	V	V
	Parallel Programming		Stream Processing(Timestream)	Processing Pipeline (SEDA)	X	V	X	V	V
	Maintainability		Improve-meents	Algorithmic Skeletons, SOA, Microservices Dynamic Distribution PGAS	X	V	X	V	V
Section 3.1.1		Section 3.1.2		Section 3.2.1		Section 3.2.2			

Table 3.5 – Synthesis of the state of the art in software design

3.3 Seamless Development



The section 3.1 shows that modularity and functional programming, especially higher-order programming and lazy evaluation, are the best organization to improve maintainability of an application. This organization is best supported by a functional approach. Indeed, higher-order programming improves readability and maintainability. However, higher-order programming, and modular programming in general, requires the use of a global memory store.

The section 3.2 shows that to attain scalability, an application needs to be organized to distribute its memory store into independent silos to provide isolation and immutability. Leading to multiply the exclusive accesses. Still, many works provide this global memory store interface to developers, because it is the best way to support the modularity advocated in section ???. This incompatibility between these two organization, and their goals is responsible for the shifts operated during the life of an application. Huge developing efforts are made to translate manually from one organization into the other, and to maintain the implementation despite its unmaintainable nature.

In section 3.3, we show different tentatives to reconciles the two organizations. Most are satisfactory for specific domains, such as the high-performance computing. It is profitable, as the expected speedup of developing an application with an adapted programming model compensates the huge development effort. However, none are satisfactory in the case of web applications because the need for performance is always uncertain. The development effort is not required at the beginning, hence its cost

	Maintainability	Performance	Both
General	Functional Programming	Message-passing	Loop parallelization
Web	Javascript	Pipeline architecture	\emptyset

Table 3.6 – Summary of the state of the art

cannot be justified. It is only when the audience increases, often with the revenue, that the cost for the development effort can be justified. This situation illustrate the need for a programming model reconciling the two concerns, of maintainability and performance.

Our objectives is to provide seamless development. 3.1 shows that modularity is not scalable but needed. But compilation can help to get scalable. 3.2 shows that parallelism is not maintainable (no Higher-Order Programming, no lazy evaluation, so no good glue between modules) But good developers can maintain double representation. With help from an equivalence or a *compiler* most developers could develop using this double representation, and seamlessly achieving parallelism and maintainability.

* Our objectives is to find an equivalence between these two organization, specifically for the case of web applications. To do so, we focus on the Javascript programming language, and specifically, the node.js interpreter. * As explained in the end of chapter 2, the execution model of Javascript is similar to a pipeline. We intend to split a node.js application into a parallel pipeline of stages.

The contribution of this thesis is organized in two chapters, as illustrated in figure 3.1. In chapter 4 I present the extraction of a pipeline of operations from a Javascript application. I show that such pipeline is similar to the one exposed by Promises, and I propose a simpler alternative to the latter called Dues. However, these operations still require a global memory for coordination so they are not executed in parallel. In chapter 5, we present the isolation of the operations into isolated containers called Fluxions.

————— need integration —————

We show that there is no languages that features higher-order functions to improve modularity, a common memory store easy to develop with, but at the same time provides scalable concurrency.

We aim at filling this gap, and for a concrete example, focus our work on the Javascript programing language. Indeed, Javascript features higher-order functions, is highly-used in concurrent context, but lacks scalable concurrency.



integrate
these
paragraphs
after
the
next
section



TODO
clear

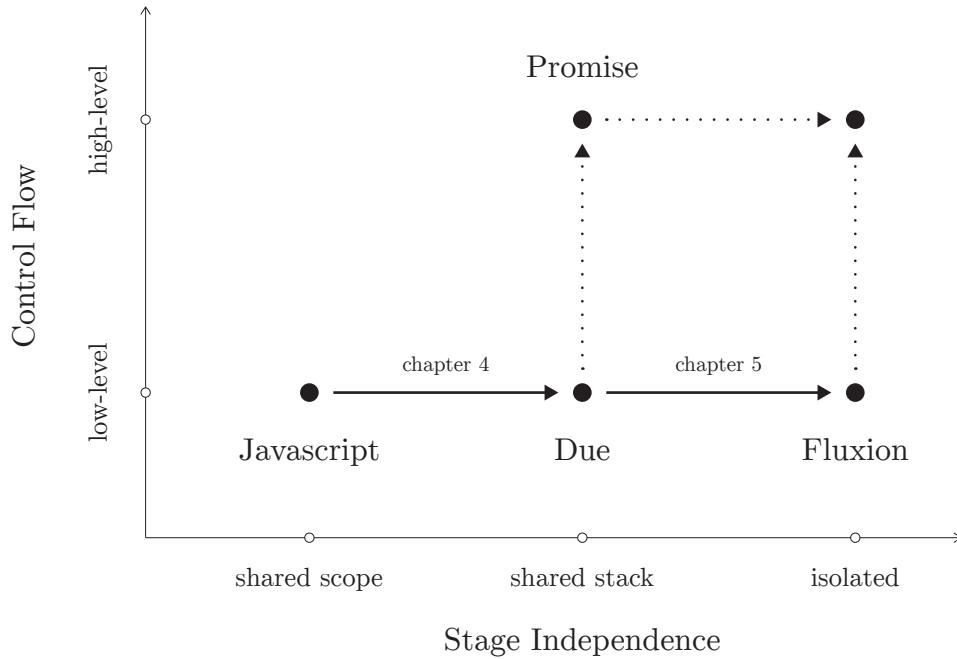


Figure 3.1 – Roadmap for this work

3.3.1 Equivalence

3.3.1.1 Rupture Point

The execution of the pipeline architecture is well delimited in isolated stages. Each stage has its own thread of execution, and is independent from the others. On the contrary, the code of the event-loop is linear because of the continuation passing style and the common memory store. However, the execution of the different callbacks are as distinct as the execution of the different stages of a pipeline. The call stacks of two callbacks are distincts. Therefore, an asynchronous function call represents the rupture between two call stacks. It is a rupture point, and is equivalent to a data stream between two stages in the pipeline architecture.

Both the pipeline architecture and the event-loop present these rupture points. The detection of rupture points allows to map a pipeline architecture onto the implementation following the event-loop model. To allow the transformation from one to the other, this thesis studies the possibility to detect rupture points, and to distribute the global memory into the parts defined by these rupture points. The detection of rupture points is addressed in chapter 4.

3.3.1.2 Invariance

The transformation should preserve the invariance as expressed by the developer to assure the correctness of the execution. The partial ordering of events in a system, by opposition to total ordering, is sufficient to assure this correctness. The global memory is a way to assure the total ordering of events, and the message passing coordination is a way to assure partial ordering of events. Therefore, to assure the correctness of the execution of a system, the state coordination with a global memory is equivalent to message passing coordination. And it is possible, at least for some rupture points, to transform the global memory coordination into message passing while conserving the correctness of execution.

In order to preserve the invariance assured by the event-loop model after the transformation, each stage of the pipeline needs to have an exclusive access to memory. The global memory needs not to be split into parts and distributed into each of the stages. To assure the missing coordinations assured by the shared memory between the stages, the transformation should provide equivalent coordination with message passing. The isolation and replacement of the global memory is fully address in chapter 5

Chapter 4

Pipeline extraction

The previous chapter presented globally the state of the art in designing systems to scale in performance, and in maintenance. It refined the scope of this thesis to the study of the opposition between maintenance scalability and performance scalability in streaming web applications. It concluded with the objectives of this thesis, which is to find an equivalence between the two opposed organizations. The maintenance scalability organization, supported by modular programming, higher-order programming and a global memory store. The performance scalability organization, supported by the parallelism of memory and exution distribution. The equivalence between these two organization is in two steps, as presented in figure 3.1. This chapter presents the first step in this equivalence. That is to identify and extract a pipeline of execution inside an application following the first organization. In this work, we focus on Javascript, and specifically node.js applications. In this chapter, I define further the higher-order programming concepts.

In Javascript, functions are first-class citizens ; it allows to manipulate them like any object, and to link them to react to asynchronous events, *e.g.* user inputs and remote requests. These asynchronously triggered functions are named callbacks, and allow to efficiently cope with the distributed and inherently asynchronous architecture of the Internet. To execute a suite of asynchronous functions, callbacks are nested one into the other. This nesting, if not organized properly, can result in unreadable layer of callbacks, commonly presented as *callback hell*¹, or *pyramid of doom*.

Promises are another way to organize a suite of asynchronous operations avoiding this callback hell. They organize the operations as a well-defined pipeline. Moreover, Promises provide additional control over the asynchronous execution flow, than call-

¹<http://maxogden.github.io/callback-hell/>

backs. They are part of the next version of the Javascript language, ECMAScript 6². To avoid the equivalence being unnecessarily incomplete, we present an alternative to Promise, called Due. Due organize the operations like Promises, as a well-defined pipeline, while discarding the unnecessary additional control over the asynchronous flow.

This chapter present an equivalence, and a compiler to identify the pipeline of operating underlying in a Javascript application using callbacks, and extract it to express it as Dues. This compiler has been tested over 64 *Node.js* packages from the node package manager (npm³). 55 packages were incompatible with the compiler, 9 packages were compiled with success.

Callbacks, Promises and Dues are further defined in section 4.1. Section 4.2 explains the transformation from imbrications of callbacks to sequences of Dues. Section 5.2 presents a compiler to automate the application of this equivalence. And finally, the developed compiler is evaluated in section 5.3.

4.1 Definitions

4.1.1 Callback

A callback is a function passed as a parameter to a function call. It is invoked by the callee to continue the execution with data not available in the caller context. We distinguish three kinds of callbacks.

Iterators are functions called for each item in a set, often synchronously.

Listeners are functions called asynchronously for each event in a stream.

Continuations are functions called asynchronously once a result is available.

As we will see later, Promises are designed as placeholders for a unique outcome. Iterators and Listeners are invoked multiple times resulting in multiple outcomes. Only continuations are equivalent to Promises. Therefore, we focus on continuations in this paper.

Callbacks are often mistaken for continuations; the former are not inherently asynchronous while the latter are. In a synchronous paradigm, the sequentiality of the execution flow is trivial. An operation needs to complete before executing the

²<http://people.mozilla.org/~jorendorff/es6-draft.html>

³<https://www.npmjs.com/>

next one. In an asynchronous paradigm, parallelism is trivial, but the sequentiality of operations needs to be explicit. Continuations are the functional way of providing this control over the sequentiality of the asynchronous execution flow.

A continuation is a function passed as an argument to allow the callee not to block the caller until its completion. The caller is able to continue the execution while the callee runs in background. The continuation is invoked later, at the termination of the callee to continue the execution as soon as possible and process the result; hence the name continuation. It provides a necessary control over the asynchronous execution flow. It also brings a control over the data flow which essentially replaces the `return` statement at the end of a synchronous function. At its invocation, the continuation retrieves both the execution flow and the result.

The convention on how to hand back the result must be common for both the callee and the continuation. For example, in *Node.js*, the signature of a continuation uses the *error-first* convention. The first argument contains an error or `null` if no error occurred; then follows the result. Listing 4.1 is a pattern of such a continuation. However, continuations don't impose any conventions; indeed, other conventions are used in the browser.

```
1 my_fn(input, function continuation(error, result) {  
2   if (!error) {  
3     console.log(result);  
4   } else {  
5     throw error;  
6   }  
7 });
```

Listing 4.1 – Example of a continuation

The callback hell occurs when many asynchronous calls are arranged to be executed sequentially. Each consecutive operation adds an indentation level, because it is nested inside the continuation of the previous operation. It produces an imbrication of calls and function definitions, as shown in listing 4.2. We say that continuations lack the chained composition of multiple asynchronous operations. Promises allow to arrange such a sequence of asynchronous operations in a more concise and readable way.

```
1 my_fn_1(input, function cont(error, result) {  
2   if (!error) {  
3     my_fn_2(result, function cont(error, result) {  
4       if (!error) {  
5         my_fn_3(result, function cont(error, result) {  
6           if (!error) {  
7             console.log(result);  
8           } else {  
9             throw error;  
10            }  
11          }  
12        }  
13      }  
14    }  
15  }  
16});
```

```

12     } else {
13         throw error;
14     }
15 });
16 } else {
17     throw error;
18 }
19 );

```

Listing 4.2 – Example of a sequence of continuations

4.1.2 Promise

In a synchronous paradigm, the sequentiality of the execution flow is trivial. While in an asynchronous paradigm, this control is provided by continuations. Promises provide a unified control over the execution flow for both paradigms. The ECMAScript 6 specification⁴ defines a Promise as an object that is used as a placeholder for the eventual outcome of a deferred (and possibly asynchronous) operation. Promises expose a `then` method which expects a continuation to continue with the result; this result being synchronously or asynchronously available.

Promises force another control over the execution flow. According to the outcome of the operation, they call one function to continue the execution with the result, or another to handle errors. This conditional execution is indivisible from the Promise structure. As a result, Promises impose a convention on how to hand back the outcome of the deferred computation, while classic continuations leave this conditional execution to the developer.

```

1 var promise = my_fn_pr(input)
2
3 promise.then(function onSuccess(result) {
4     console.log(result);
5 }, function onError(error) {
6     throw error;
7 });

```

Listing 4.3 – Example of a promise

Promises are designed to fill the lack of chained composition from continuations. They allow to arrange successions of asynchronous operations as a chain of continuations, by opposition to the imbrication of continuations illustrated in listing 4.2. That is to arrange them, one operation after the other, in the same indentation level.

The listing 4.4 illustrates this chained composition. The functions `my_fn_pr_2` and `my_fn_pr_3` return promises when they are executed, asynchronously. Because

⁴<https://people.mozilla.org/~jorendorff/es6-draft.html#sec-promise-objects>

these promises are not available synchronously, the method `then` synchronously returns intermediary Promises. The latter resolve only when the former resolve. This behavior allows to arrange the continuations as a flat chain of calls, instead of an imbrication of continuations.

```

1 my_fn_pr_1(input)
2 .then(my_fn_pr_2, onError)
3 .then(my_fn_pr_3, onError)
4 .then(console.log, onError);
5
6 function onError(error) {
7   throw error;
8 }
```

Listing 4.4 – A chain of Promises is more concise than an imbrication of continuations

The Promises syntax is more concise, and also more readable because it is closer to the familiar synchronous paradigm. Indeed, Promises allow to arrange both the synchronous and asynchronous execution flow with the same syntax. It allows to easily arrange the execution flow in parallel or in sequence according to the required causality. This control over the execution leads to a modification of the control over the data flow. Programmers are encouraged to arrange the computation as series of coarse-grained steps to carry over inputs. In this sense, Promises are comparable to some coarse-grained data-flow programming paradigms, such as Flow-based programming [58].

4.1.3 From continuations to Promises

As detailed in the previous sections, continuations provide the control over the sequentiality of the asynchronous execution flow. Promises improve this control to allow chained compositions, and unify the syntax for the synchronous and asynchronous paradigm. This chained composition brings a greater clarity and expressiveness to source codes. At the light of these insights, it makes sense for a developer to switch from continuations to Promises. However, the refactoring of existing code bases might be an operation impossible to carry manually within reasonable time. We want to automatically transform an imbrication of continuations into a chained composition of Promises.

We identify two steps in this transformation. The first is to provide an equivalence between a continuation and a Promise. The second is the composition of this equivalence. Both steps are required to transform imbrications of continuations into chains of Promises.

Because Promises bring chained composition, the first step might seem trivial as it does not imply any imbrication to transform into chain. However, as explained in

section 4.1.2, Promises impose a control over the execution flow that continuations leave free. This control induces a common convention to hand back the outcome to the continuation.

In the Javascript landscape, there is no dominant convention for handing back outcomes to continuations. In the browser, many conventions coexist. For example, *jQuery*'s `ajax`⁵ method expects an object with different continuations for success, errors and various other events during the asynchronous operation. *Q*⁶, a popular library to control the asynchronous flow, exposes two methods to define continuations: `then` for successes, and `catch` for errors. On the other hand, the *Node.js* API always used the *error-first* convention, encouraging developers to provide libraries using the same convention. In this large ecosystem the *error-first* convention is predominant. All these examples use different conventions than the Promise specification detailed in section 4.1.2. They present strong semantic differences, despite small syntactic differences.

To translate these different conventions into the Promises one, the compiler would need to identify them. Such an identification might be possible with static analysis methods such as the points-to analysis [77], or a program logic [27, 13]. However, it seems impracticable because of the number and semantical heterogeneity of these conventions. Indeed, in the browser, each library seems to provide its own convention.

In this paper, we are interested in the transformation from imbrications to chains, not from one convention to another. The *error-first* convention, used in *Node.js*, is likely to represent a large, coherent code base to test the equivalence. Indeed contains currently more than 125 000 packages. For this reason, we focus only on the *error-first* convention. Thus, our compiler is only able to compile code that follows this convention. The convention used by Promises is incompatible. We propose an alternative specification to Promise following the *error-first* convention. In the next section we present this specification called Due.

The choice to focus on *Node.js* is also motivated by our intention to compare later the chained sequentiality of Promises with the data-flow paradigm. *Node.js* allows to manipulate streams of messages. This proved to be efficient for real-time web applications manipulating streams of user requests. Both Promises and data-flow arrange the computation in chains of independent operations.

⁵<http://api.jquery.com/jquery.ajax/>

⁶<http://documentup.com/kriskowal/q/>

4.1.4 Due

A Due is an object used as placeholder for the eventual outcome of a deferred operation. Dues are a simplification of the Promise specification. They are essentially similar to Promises, except for the convention to hand back outcomes. They use the *error-first* convention, like *Node.js*, as illustrated in listing 4.5. The implementation of Dues and its tests are available online⁷. A more in-depth description of Dues and their creation follows in the next paragraphs.

```
1 var my_fn_due = require('due').mock(my_fn);
2
3 var due = my_fn_due(input);
4
5 due.then(function continuation(error, result) {
6   if (!error) {
7     console.log(result);
8   } else {
9     throw error;
10 }
11});
```

Listing 4.5 – Example of a due

A due is typically created inside the function which returns it. In listing 4.5, line 1, the `mock` method wraps `my_fn` in a Due-compatible function. The rest of this code is similar to the Promise example, listing 4.3.

We illustrate in listing 4.6 the creation of a Due through the `mock` method. At its creation, line 6, the Due expects a callback containing the deferred operation, which is `my_fn` here. This callback is executed synchronously with the function `settle` as argument to settle the Due, synchronously or asynchronously. The `settle` function is pushed at the end of the list of arguments. The callback invokes the deferred operation with this list of arguments, and the current context, line 8. When finished, the latter calls `settle` to settle the Due and save the outcome. Settled or not, the created Due is always synchronously returned. Its `then` method allows to define a continuation to retrieve the saved outcome, and continue the execution after its settlement. If the deferred operation is synchronous, the Due settles during its creation and the `then` method immediately calls this continuation. If the deferred operation is asynchronous, this continuation is called during the Due settlement.

```
1 Due.mock = function(my_fn) {
2   return function mocked_fn() {
3     var _args = Array.prototype.slice.call(arguments),
4       _this = this;
5
6     return new Due(function(settle) {
7       _args.push(settle);
```

⁷<https://www.npmjs.com/package/due>

```

8     my_fn.apply(_this, _args);
9   })
10 }
11 }
```

Listing 4.6 – Creation of a due

The composition of Dues is the same than for Promises (see section 4.1.2). Through this chained composition, Dues arrange the execution flow as a sequence of actions to carry on inputs.

This simplified specification adopts the same convention than *Node.js* for continuations to hand back outcomes. Therefore, the equivalence between a continuation and a Due is trivial. Dues are admittedly tailored for this paper, hence, they are not designed to be written by developers, like Promises are. They are an intermediary step between classical continuations and Promises. We present in section 4.2 the equivalence between continuations and Dues.

4.2 Equivalence

*



TODO this title is not clear

We present the transformation from a nested imbrication of continuations into a chain of Dues. We explain the three limitations imposed by our compiler for this transformation to preserve the semantic. They preserve the execution order, the execution linearity and the scopes of the variables used in the operations.

4.2.1 Execution order

Our compiler spots function calls with a continuation, which are similar to the abstraction in (4.1). It wraps the function *fn* into the function *fn_{due}* to return a Due. And it relocates the continuation in a call to the method **then**, which references the Due previously returned. The result should be similar to (4.2). The differences are highlighted in bold font.

$$fn([arguments], continuation) \tag{4.1}$$

$$fn_{\mathbf{due}}([arguments]).\mathbf{then}(continuation) \tag{4.2}$$

The execution order is different whether *continuation* is called synchronously, or asynchronously. If *fn* is synchronous, it calls the *continuation* within its execution. It might execute *statements* after executing *continuation*, before returning. If *fn* is

asynchronous, the continuation is called after the end of the current execution, after fn . The transformation erases this difference in the execution order. In both cases, the transformation relocates the execution of *continuation* after the execution of fn . For synchronous fn , the execution order changes ; the execution of *statements* at the end of fn and the continuation switch. The latter must be asynchronous to preserve the execution order.

4.2.2 Execution linearity

Our compiler transforms a nested imbrication of continuations, which is similar to the abstraction in (4.3) into a flatten chain of calls encapsulating them, like in (4.4).

```

 $fn1([arguments], cont1\{
    declare variable \leftarrow result
    fn2([arguments], cont2\{
        print variable
    })
})$  (4.3)

```

```

declare variable
 $fn1_{\text{due}}([arguments])$ 
.then( $cont1\{
    variable \leftarrow result
    fn2_{\text{due}}([arguments])
\}$ )
.then( $cont2\{
    print variable
\}$ )
```

An imbrication of continuations must not contain any loop, nor function definition that is not a continuation. Both modify the linearity of the execution flow which is required for the equivalence to keep the semantic. A call nested inside a loop returns multiple Dues, while only one is returned to continue the chain. A function definition breaks the execution linearity. It prevent the nested call to return the Due expected to continue the chain. On the other hand, conditional branching leaves the execution linearity and the semantic intact. If the nested asynchronous function is not called, the execution of the chain stops as expected.

4.2.3 Variable scope

In (4.3), the definitions of *cont1* and *cont2* are overlapping. The *variable* declared in *cont1* is accessible in *cont2* to be printed. In (4.4), however, definitions of *cont1* and *cont2* are not overlapping, they are siblings. The *variable* is not accessible to *cont2*. It must be relocated in a parent function to be accessible by both *cont1* and *cont2*. To detect such variables, the compiler must infer their scope statically. Languages with a lexical scope define the scope of a variable statically. Most imperative languages present a lexical scope, like C/C++, Python, Ruby or Java. The subset of Javascript excluding the built-in functions `with` and `eval` is also lexically scoped. To compile Javascript, the compiler must exclude programs using these two statements.

4.3 Compiler

We build a compiler to automate the application of this equivalence on existing Javascript projects. The compilation process contains two important steps, the identification of the continuations, and the generation of chains.

4.3.1 Identification of continuations

The first compilation step is to identify the continuations and their imbrications. The nested imbrication of callbacks only occurs when they are defined *in situ*. The compiler detects a function definition within the arguments of a function call. This detection is based on the syntax, and is trivial.

Not all detected callbacks are continuations, but the equivalence is applicable only on the latter. A continuation is a callback invoked only once, asynchronously. Spotting a continuation implies to identify these two conditions. There is no syntactical difference between a synchronous and an asynchronous callee. And it is impossible to assure a callback to be invoked only once, because the implementation of the callee is often statically unavailable. Therefore, the identification of continuations is necessarily based on semantical differences. To recognize these differences, the compiler would need to have a deep understanding of the control and data flows of the program. Because of the highly dynamic nature of Javascript, this understanding is either unsound, limited, or complex. Instead, we choose to leave to the developer the identification of compatible continuations among the identified callbacks. They are expected to understand the limitations of this compiler, and the semantic of the code to compile.

We provide a simple interface for developers to interact with the compiler. We built this interface around the compiler in a web page available online⁸ to reproduce the tests. The web technologies allow to quickly build an interface for a wide variety of computing devices.

This interaction prevents the complete automation of the individual compilation process. However, we are working on an automation at a global scale. We expect to be able to identify a continuation only based on the name of its callee, *e.g.* `fs.readFile`. We built a service to gather these names along with their identification. The compiler queries this service to present to the developer an estimated identification. After the compilation, it sends back the identification corrected by the developer to refine the future estimations. In future works, we would like to study the possibility for such a service to assist, and ease the compilation process.

4.3.2 Generation of chains

The compositions of continuations and Dues are arranged differently. Continuations structure the execution flow as a tree, while a chain of Dues imposes to arrange it sequentially. A parent continuation can execute several children, while a Due allow to chain only one. The second compilation step is to identify the imbrications of continuations, and trim the extra branches to transform them into chains.

If a continuation has more than one child, the compiler tries to find a single legitimate child to form the longest chain possible. This legitimate child is the only parent among its siblings. If there are several parents among the children, none are the legitimate child. The non legitimate children start a new tree. This step transform each tree of continuations into several chains of continuations that translate into sequences of Dues. The code generation from these chains is straightforward from the equivalence.

4.4 Evaluation

To validate our compiler, we compile several Javascript projects likely to contain continuations. We present the results of these tests.

The compilation of a project requires user interaction. To conduct the test in a reasonable time, we limit the test set to a minimum. We search the *Node Package Manager* database to restrict the set to *Node.js* projects. We refine the selection to web applications depending on the web framework *express*, but not on the most

⁸compiler-due.apps.zone52.org

common Promises libraries such as *Q* and *Async*. We refine further the selection to projects using the test frameworks *mocha* in its default configuration. We use these tests to validate the compiler. The test set contains 64 projects. This subset is very small, and cannot represent the wide possibilities of Javascript. However, we believe it is sufficient to represent a majority of common cases.

For each project, we verify that it is correctly tested, and passes the tests. During the compilation, we identify the compatible continuations among the detected callbacks. We apply the unmodified test on the compilation result. The compilation result should pass the tests as well. This is not a strong validation, but it assures the compiler to work as expected in most common cases.

Of the 64 projects tested, almost a half, does not contain any compatible continuations. We reckon that these projects use continuations the compiler is unable to detect. The other projects were rejected by the compiler because they contain `with` or `eval` statements, they use Promises libraries we didn't filter previously. 9 projects compiled successfully. The compiler did not fail to compile any project of the initial test set.

Over the 9 successfully compiled projects, the compiler detected 172 callbacks. We manually identified 56 of them to be compatible continuations. The false positives are mainly the listeners that the web applications register to react to user requests.

One project contains 20 continuations, the others contains between 1 and 9 continuations each. On the 56 continuations, 36 are single. The others 20 continuations belong to imbrications of 2 to 4 continuations. The result of this evaluation prove the compiler to be able to successfully transform imbrications of continuations.

On the 64 projects composing the test set

29 (45.3%) do not contain any compatible continuations,

10 (15.6%) are not compilable because they contain `with` or `eval` statements,

5 (7.8%) use less common asynchronous libraries we didn't filter previously,

4 (6.3%) are not syntactically correct,

4 (6.3%) fail their tests before the compilation,

3 (4.7%) are not tested, and

10 (14.0%) compile successfully.

The compiler do not fail to compile any project. The details of these projects are available in Appendix ??.

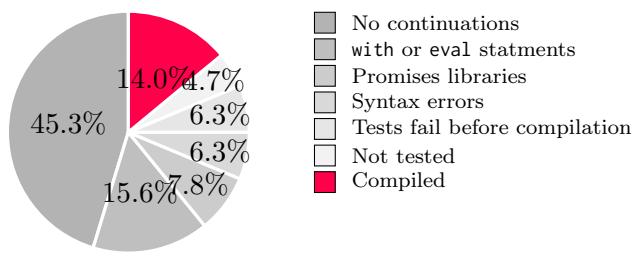


Figure 4.1 – Compilation results distribution

Chapter 5

Pipeline isolation

The previous chapter presented a compiler to identify and extract the underlying pipeline in a Javascript application. However, all the operations are not independent, and cannot be executed in parallel, to support the performance scalability. This chapter present the second contribution of this thesis. The equivalence between a memory shared among all the operations and independent memory for each operation in a pipeline. It tackles the problems arising from the translation of the global memory synchronization into message passing.

This equivalence is implemented as a compiler, improving upon the previous one. The compiler transforms a Javascript application into a network of independent parts communicating by message streams and executed in parallel. We named these parts *fluxions*, by contraction between a flux and a function.

Section 5.1 describes the execution model that executes fluxions in parallel, and assure their communications. The compiler, and the equivalence are described in section 5.2. Section 5.3 a real-case test of compilation, and expose the limits of this compiler.

5.1 Fluxional execution model

In this section, we present an execution model to provide scalability to web applications. To achieve this, the execution model provides a granularity of parallelism at the function level. Functions are encapsulated in autonomous execution containers with their state, so as to be reallocated and executed in parallel. This execution model is close to the actors model, as the execution containers are independent and communicate by messages. The communications are assimilated to stream of messages, similarly to the dataflow programming model. It allows to reason on the

throughput of these streams, and to react to load increases.

The fluxional execution model executes programs written in our high-level fluxional language, whose grammar is presented in figure 5.1. An application \langle program \rangle is partitioned into parts encapsulated in autonomous execution containers named *fluxions* \langle flx \rangle . In the following paragraphs, we present the *fluxions*. Then we present the messaging system to carry the communications between *fluxions*. Finally, we present an example application using this execution model.

5.1.1 Fluxions

A *fluxion* \langle flx \rangle is named by a unique identifier \langle id \rangle to receive messages, and might be part of one or more groups indicated by tags \langle tags \rangle . A *fluxion* is composed of a processing function \langle fn \rangle , and a local memory called a *context* \langle ctx \rangle . At a message reception, the *fluxion* modifies its *context*, and sends messages on its output streams \langle streams \rangle to downstream *fluxions*. The *context* handles the state on which a *fluxion* relies between two message receptions. In addition to message passing, the execution model allows *fluxions* to communicate by sharing state between their *contexts*. The fluxions that need to synchronize together are grouped with the same tag, and loose their independence.

There are two types of streams, *start* and *post*, which correspond to the nature of the rupture point yielding the stream. We differentiate the two types with two different arrows, double arrow ($>>$) for *start* rupture points and simple arrow ($->$) for *post* rupture points. The two types of rupture points are further detailed in section 5.2.1.1.

5.1.2 Messaging system

The messaging system assures the stream communications between fluxions. It carries messages based on the names of the recipient fluxions. After the execution of a fluxion, it queues the resulting messages for the event loop to process.

The execution cycle of an example fluxional application is illustrated in figure 5.2. Circles represent registered fluxions. The source code for this application is in listing 5.1 and the fluxional code for this application is in listing 5.2. The fluxion *reply* has a context containing the variable `count` and `template`. The plain arrows represent the actual message paths in the messaging system, while the dashed arrows between fluxions represent the message streams as seen in the fluxionnal application.

The *main* fluxion is the first fluxion in the flow. When the application receives a request, this fluxion triggers the flow with a `start` message containing the request,

$$\begin{aligned}
\langle \text{program} \rangle &\equiv \langle \text{flx} \rangle \mid \langle \text{flx} \rangle \text{ eol} \langle \text{program} \rangle \\
\langle \text{flx} \rangle &\equiv \text{f1x} \langle \text{id} \rangle \langle \text{tags} \rangle \langle \text{ctx} \rangle \text{ eol} \langle \text{streams} \rangle \text{ eol} \langle \text{fn} \rangle \\
\langle \text{tags} \rangle &\equiv \& \langle \text{list} \rangle \mid \text{empty string} \\
\langle \text{streams} \rangle &\equiv \text{null} \mid \langle \text{stream} \rangle \mid \langle \text{stream} \rangle \text{ eol} \langle \text{streams} \rangle \\
\langle \text{stream} \rangle &\equiv \langle \text{type} \rangle \langle \text{dest} \rangle [\langle \text{msg} \rangle] \\
\langle \text{dest} \rangle &\equiv \langle \text{list} \rangle \\
\langle \text{ctx} \rangle &\equiv \{ \langle \text{list} \rangle \} \\
\langle \text{msg} \rangle &\equiv [\langle \text{list} \rangle] \\
\langle \text{list} \rangle &\equiv \langle \text{id} \rangle \mid \langle \text{id} \rangle , \langle \text{list} \rangle \\
\langle \text{type} \rangle &\equiv \text{>>} \mid \text{->} \\
\langle \text{id} \rangle &\equiv \text{Identifier} \\
\langle \text{fn} \rangle &\equiv \text{imperative language and stream syntax}
\end{aligned}$$

Figure 5.1 – Syntax of a high-level language to represent a program in the fluxionnal form

②. This first message is to be received by the next fluxion handler, ③ and ④. The fluxion handler sends back a message, ⑤, to be enqueued, ⑥. The system loops through steps ③ through ⑥ until the queue is empty. This cycle starts again for each new incoming request causing another `start` message.

5.1.3 Service example

To illustrate the fluxional execution model, and the compiler, we present in listing 5.1 an example of a simple web application. This application reads a file, and sends it back along with a request counter.

```

1 var app = require('express')(),
2   fs = require('fs'),
3   count = 0;
4
5 app.get('/', function handler(req, res){
6   fs.readFile(__filename, function reply(err, data) {
7     count += 1;
8     res.send(err || template(count, data));
9   });
10 });
11
12 app.listen(8080);

```

Listing 5.1 – Example web application

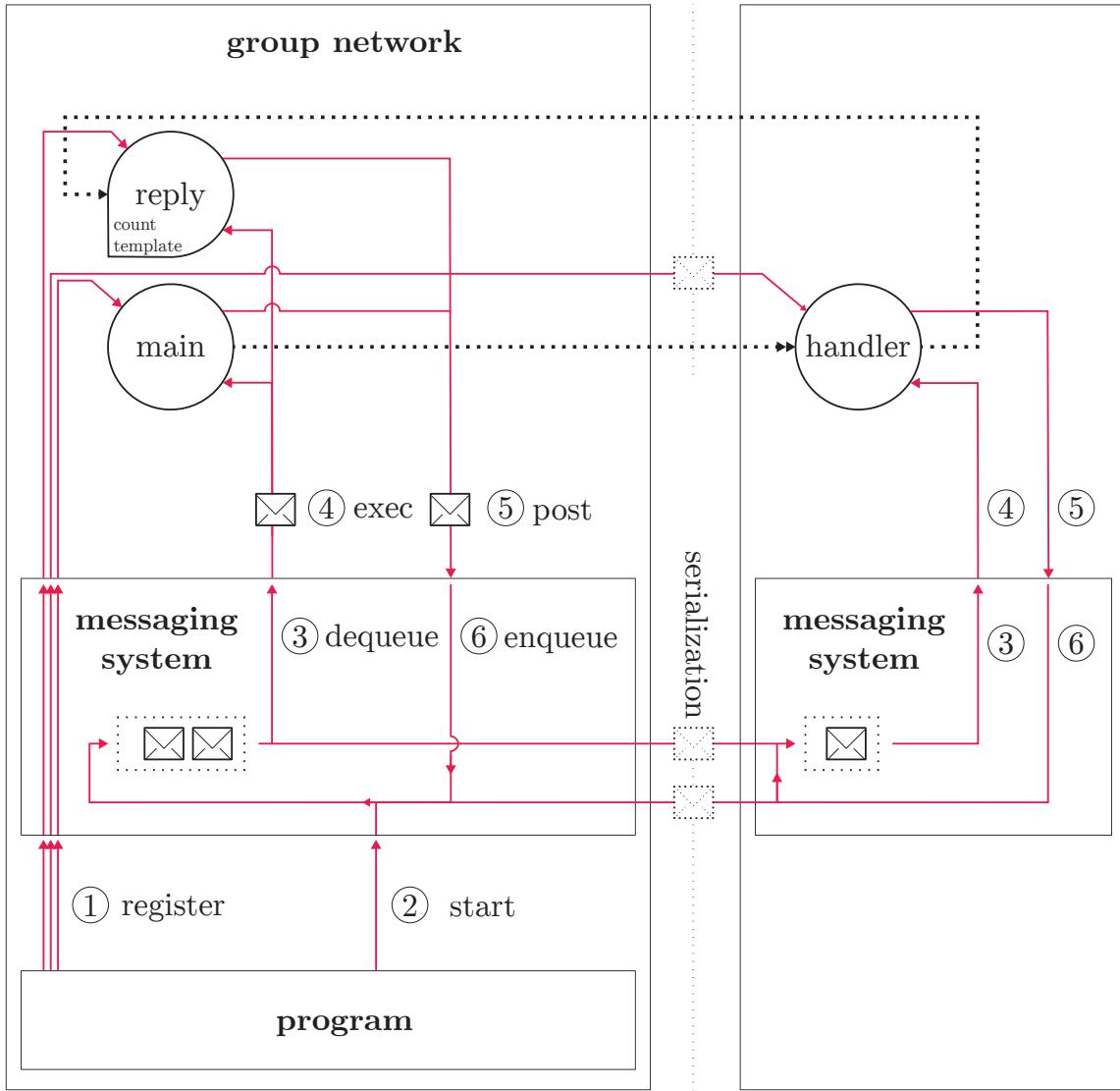


Figure 5.2 – The fluxionnal execution model in details

The **handler** function, line 5 to 11, receives the input stream of request. The **count** variable at line 3 increments the request counter. This object needs to be persisted in the fluxion *context*. The **template** function formats the output stream to be sent back to the client. The **app.get** and **res.send** functions, respectively line 5 and 8, interface the application with the clients. And between these two interface functions is a chain of three functions to process the client requests : **app.get** →

→ **handler** → **reply**. This application is transformed into the high-level fluxionnal language in listing 5.2 which is illustred in Figure 5.2.

```

1 flx main & network
2 >> handler [res]
3   var app = require('express')(),
4     fs = require('fs'),
5     count = 0;
6
7   app.get('/', >> handler); //
8   app.listen(8080);
9
10 flx handler
11 -> reply [res]
12   function handler(req, res) {
13     fs.readFile(__filename, -> reply); //
14   }
15
16 flx reply & network {count, template}
17 -> null
18   function reply(error, data) {
19     count += 1; //
20     res.send(err || template(count, data)); //
21 }
```

Listing 5.2 – Example application expressed in the high-level fluxional language

The application is organized as follow. The flow of requests is received from the clients by the fluxion **main**, it continues in the fluxion **handler**, and finally goes through the fluxion **reply** to be sent back to the clients. The fluxions **main** and **reply** have the tag **network**. This tag indicates their dependency over the network interface, because they received the response from and send it back to the clients. The fluxion **handler** doesn't have any dependencies, hence it can be executed in parallel.

The last fluxion, **reply**, depends on its context to holds the variable **count** and the function **template**. It also depends on the variable **res** created by the first fluxion, **main**. This variable is carried by the stream through the chain of fluxion to the fluxion **reply** that depends on it. This variable holds the references to the network sockets. It is the variable the group **network** depends on.

Moreover, if the last fluxion, **reply**, did not relied on the variable **count**, the group **network** would be stateless. The whole group could be replicated as many time as needed.

This execution model allows to parallelize the execution of an application. Some parts are arranged in pipeline, like the fluxion **handler**, some other parts are replicated, as could be the group **network**. This parallelization improves the scalability of the application. Indeed, as a fluxion contains its state and expresses its dependencies, it can be migrated. It allows to adapt the number of fluxions per core to

adjust the resource usage in function of the desired throughput.

Our goal, as described in the introduction, is not to propose a new programming paradigm with this high-level language but to automate the architecture shift. We present the compiler to automate this architecture shift in the next section.

5.2 Fluxionnal compiler

The source languages we focus on should present higher-order functions and be implemented as an event-loop with a global memory. Javascript is such a language : it doesn't require an event-loop, but it is often implemented on top of an event-loop. *Node.js* is an example of such an implementation. We developed a compiler that transforms a *Node.js* application into a fluxional application compliant with the execution model described in section 5.1.

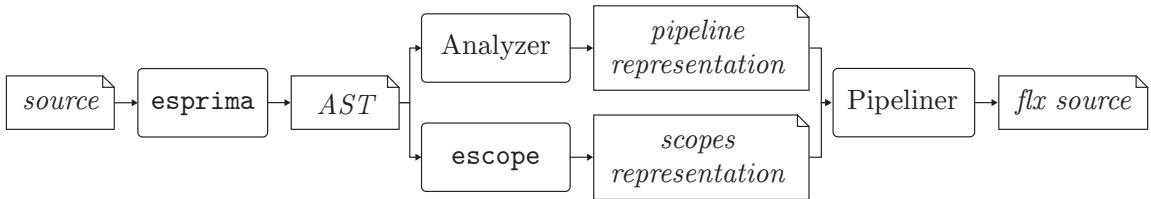


Figure 5.3 – Compilation chain

The chain of compilation is described in figure 5.3. From the source of a *Node.js* application, the compiler extracts an Abstract Syntax Tree (AST) with *esprima*. From this AST, the analyzer step identifies the limits of the different application parts and how they relate to form a pipeline. This first step outputs a pipeline representation of the application. Section 5.2.1 explains this first compilation step. In the pipeline representation, the stages are not yet independent and encapsulated into fluxions. From the AST, *escope* produces a representation of the memory scopes. The pipeliner step analyzes the pipeline representation and the scopes representation to distribute the shared memory into independent groups of fluxions. Section 5.2.2 explains this second compilation step.

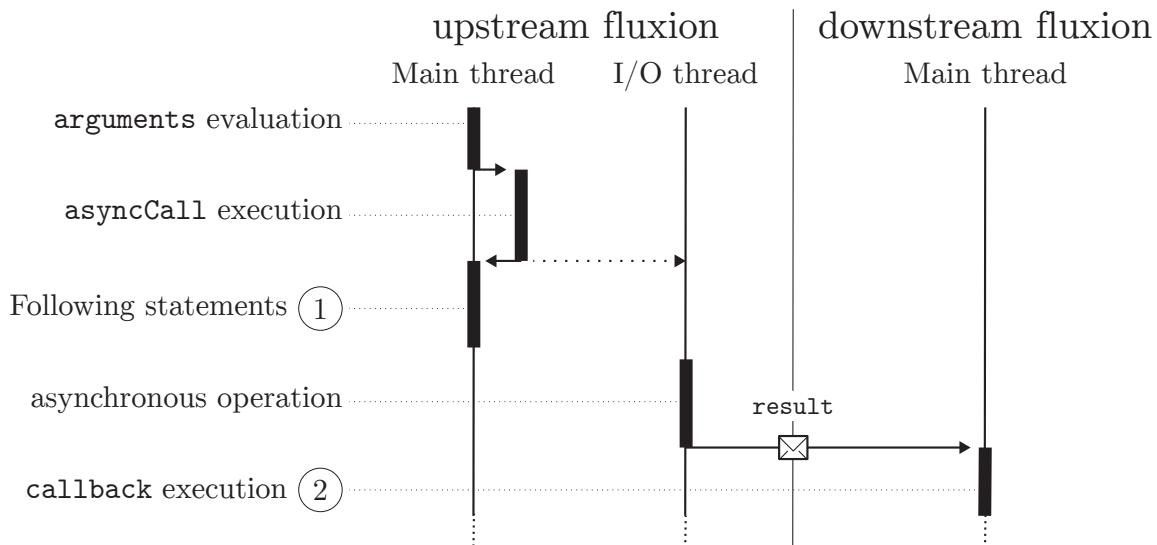
5.2.1 Analyzer step

The limit between two application parts is defined by a rupture point. The analyzer identifies these rupture points, and outputs a representation of the application in a

pipeline form, with application parts as the stages, and rupture points as the message streams of this pipeline.

5.2.1.1 Rupture points

A rupture point is a call of a loosely coupled function. It is an asynchronous call without subsequent synchronization with the caller. In *Node.js*, I/O operations are asynchronous functions and indicate such rupture point between two application parts. Figure 5.4 shows an example of a rupture point with the execution of the two application parts isolated into fluxions. The two application parts are the caller of the asynchronous function call on one hand, and the callback provided to the asynchronous function call on the other hand.



```

1 asyncCall(arguments, function callback(result){ (2) });
2 // Following statements (1)

```

Figure 5.4 – Rupture point interface

A callback is a function passed as a parameter to a function call. It is invoked by the callee to continue the execution with data not available in the caller context. We distinguish three kinds of callbacks, but only two are asynchronous : listeners and continuations. Similarly, there are two types of rupture points, respectively *start* and *post*.

Start rupture points are indicated by listeners. They are on the border between the application and the outside, continuously receiving incoming user requests. An example of a start rupture point is in listing 5.1, between the call to `app.get()`, and its listener `handler`. These rupture points indicate the input of a data stream in the program, and the beginning of a chain of fluxions to process this stream.

Post rupture points are indicated by continuations. They represent a continuity in the execution flow after an asynchronous operation yielding a unique result, such as reading a file, or querying a database. An example of a post rupture points is in listing 5.1, between the call to `fs.readFile()`, and its continuation `reply`.

5.2.1.2 Detection

The compiler uses a list of common asynchronous callees, like the `express` and file system methods. This list can be augmented to match asynchronous callees individually for any application. To identify the callee, the analyzer walks the AST to find a call expression matching this list.

After the identification of the callee, the callback needs to be identified as well to be encapsulated in the downstream fluxion. For each asynchronous call detected, the compiler test if one of the arguments is of type `function`. Some callback functions are declared *in situ*, and are trivially detected. For variable identifier, and other expressions, the analyzer tries to detect their type. To do so, the analyzer walks back the AST to track their assignations and modifications, and to determine their last value.

5.2.2 Pipeliner step

A rupture point eventually breaks the chain of scopes between the upstream and downstream fluxion. The closure in the downstream fluxion cannot access the scope in the upstream fluxion as expected. The pipeliner step replaces the need for this closure, allowing application parts to rely only on independent memory stores and message passing. It determines the distribution using the scope representation, which represents the variables' dependencies between application parts. Depending on this representation, the compiler can replace the broken closures in three different ways. We present these three alternatives with the example figure 5.5.

Scope If a variable is modified inside only one application part in the current *post* chain, then the pipeliner adds it to the context of its fluxion.

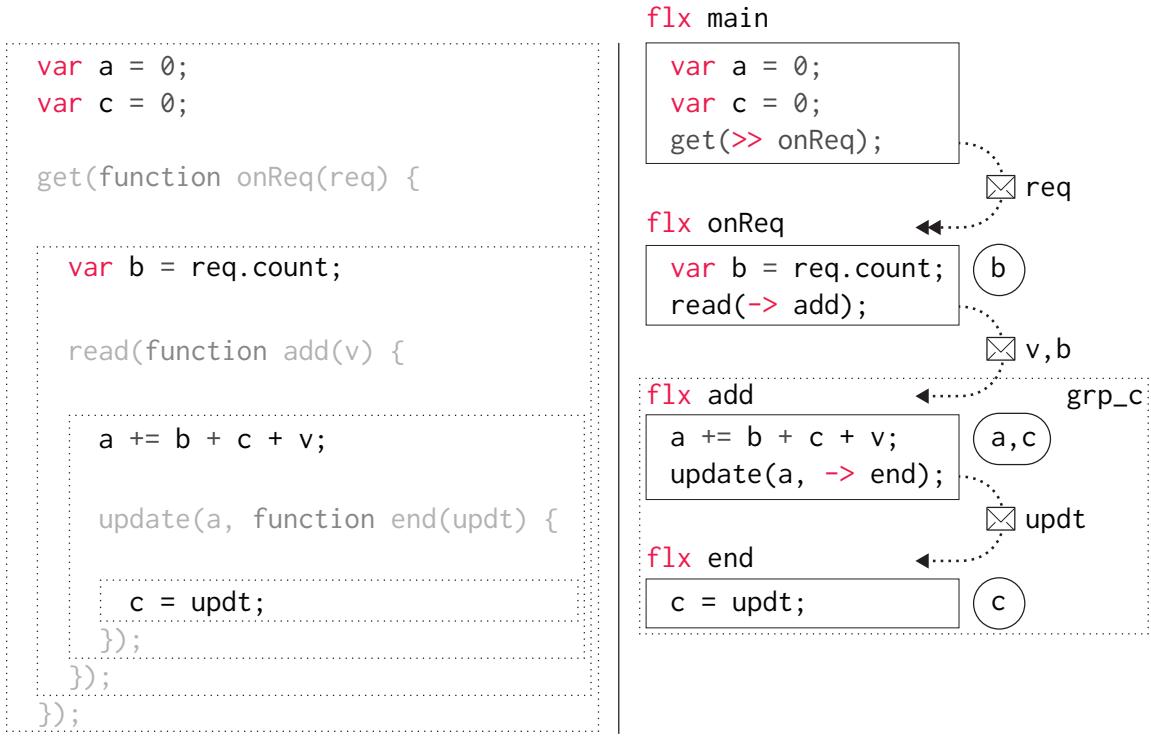


Figure 5.5 – Variable management from Javascript to the high-level fluxionnal language

In figure 5.5, the variable `a` is updated in the function `add`. The pipeliner step stores this variable in the context of the fluxion `add`.

Stream If a variable is modified inside an application part, and read inside downstream application parts, then the pipeliner makes the upstream fluxion add this variable to the message stream to be sent to the downstream fluxions. It is impossible to send variables to upstream fluxions, without race conditions. If the fluxion retro propagates the variable for an upstream fluxion to read, the upstream fluxion might use the old version while the new version is on its way.

In figure 5.5, the variable `b` is set in the function `onReq`, and read in the function `add`. The pipeliner step makes the fluxion `onReq` send the updated variable `b`, in addition to the variable `v`, in the message sent to the fluxion `add`.

Exceptionally, if a variable is defined inside a *post* chain, like `b`, then this variable can be streamed inside this *post* chain without restriction on the order of modification

and read. Indeed, the execution of the upstream fluxion for the current *post* chain is assured to end before the execution of the downstream fluxion. Therefore, no reading of the variable by the upstream fluxion happens after the modification by the downstream fluxion.

Share If a variable is needed for modification by several application parts, or is read by an upstream application part, then it needs to be synchronized between the fluxions. To respect the semantics of the source application, we cannot tolerate inconsistencies. Therefore, the pipeliner groups all the fluxions sharing this variable within a same tag. And it adds this variable to the contexts of each fluxions.

In figure 5.5, the variable `c` is set in the function `end`, and read in the function `add`. As the fluxion `add` is upstream of `end`, the pipeliner step groups the fluxion `add` and `end` with the tag `grp_c` to allow the two fluxions to share this variable.

5.3 Real case test

The goal of this test is to prove the possibility for an application to be compiled into a network of independent parts. We want to show the current limitations of this isolation and the modifications needed on the application to circumvent these limitations.

We present a test of our compiler on a real application, gifsockets-server¹. This application was selected from the `npm` registry because it depends on `express`, it is tested, working, and simple enough to illustrate this evaluation. It is part of the selection from a previous work.

This application is a real-time chat using gif-based communication channels. The server transforms the received text into a gif frame, and pushes it back to a never-ending gif to be displayed on the client. Listing 5.3 is a simplified version of this application.

```

1 var express = require('express'),
2   app = express(),
3   routes = require('gifsockets-middleware'),
4   getRawBody = require('raw-body');
5
6 function bodyParser(limit) {
7   return function saveBody(req, res, next) {
8     getRawBody(req, {
9       expected: req.headers['content-length'],
10      limit: limit
11    }, function (err, buffer) {
12      req.body = buffer;

```

¹<https://github.com/twolffson/gifsockets-server>

```

13     next();
14   });
15 }
16 }
17
18 app.post('/image/text', bodyParser(1 * 1024 * 1024), routes.writeTextToImages);
19 app.listen(8000);

```

Listing 5.3 – Simplified version of gifsockets-server

On line 18, the application registers two functions to process the requests received on the url `/image/text`. The closure `saveBody`, line 7, returned by `bodyParser`, line 6, and the method `routes.writeTextToImages` from the external module `gifsockets-middleware`, line 3. The closure `saveBody` calls the asynchronous function `getRawBody` to get the request body. Its callback handles the errors, and calls `next` to continue processing the request with the next function, `routes.writeTextToImages`.

5.3.1 Compilation

We compile this application with the compiler detailed in section 5.2. The function call `app.post`, line 18, is a rupture point. However, its callbacks, `bodyParser` and `routes.writeTextToImages` are evaluated as functions only at runtime. For this reason, the compiler ignores this rupture point, to avoid interfering with the evaluation.

The compilation result is in listing 5.4. The compiler detects a rupture point : the function `getRawBody` and its anonymous callback, line 11. It encapsulates this callback in a fluxion named `anonymous_1000`. The callback is replaced with a stream placeholder to send the message stream to this downstream fluxion. The variables `req`, and `next` are appended to this message stream, to propagate their value from the `main` fluxion to the `anonymous_1000` fluxion.

When `anonymous_1000` is not isolated from the `main` fluxion, the compilation result works as expected. The variables used in the fluxion, `req` and `next`, are still shared between the two fluxions. Our goal is to isolate the two fluxions, to be able to safely parallelize their executions.

```

1 flx main
2 >> anonymous_1000 [req, next]
3 var express = require('express'),
4     app = express(),
5     routes = require('gifsockets-middleware'), //
6     getRawBody = require('raw-body');
7
8 function bodyParser(limit) { //
9   return function saveBody(req, res, next) { //
10     getRawBody(req, { //
11       expected: req.headers['content-length'], //

```

```

12         limit: limit
13     }, >> anonymous_1000);
14   );
15 }
16
17 app.post('/image/text', bodyParser(1 * 1024 * 1024), routes.writeTextToImages); // 
18 app.listen(8000);
19
20 fix anonymous_1000
21 -> null
22 function (err, buffer) { //
23   req.body = buffer; //
24   next(); //
25 }

```

Listing 5.4 – Compilation result of gifsockets-server

5.3.2 Isolation

In listing 5.4, the fluxion `anonymous_1000` modifies the object `req`, line 23, to store the text of the received request, and it calls `next` to continue the execution, line 24. These operations produce side-effects that should propagate in the whole application, but the isolation prevents this propagation. Isolating the fluxion `anonymous_1000` produces runtime exceptions. We detail in the next paragraph, how we handle this situation to allow the application to be parallelized. This test highlights the current limitations of the compiler, and presents future works to circumvent them.

5.3.2.1 Variable `req`

The variable `req` is read in fluxion `main`, lines 10 and 11. Then it is associated in fluxion `anonymous_1000` to `buffer`, line 23. The compiler is unable to identify further usages of this variable. However, the side effect resulting from this association impacts a variable in the scope of the next callback, `routes.writeTextToImages`. We modified the application to explicitly propagate this side-effect to the next callback through the function `next`. We explain further modification of this function in the next paragraph.

5.3.2.2 Closure `next`

The function `next` is a closure provided by the `express Router` to continue the execution with the `next` function to handle the client request. Because it indirectly relies on network sockets, it is impossible to isolate its execution with the `anonymous_1000` fluxion. Instead, we modify `express`, so as to be compatible with the fluxionnal execution model. We explain the modification below.

```

1 flx main & express
2 >> anonymous_1000 [req, next]
3 var express = require('express'),
4     app = express(),
5     routes = require('gifsockets-middleware'), //
6     getRawBody = require('raw-body');
7
8 function bodyParser(limit) { //
9     return function saveBody(req, res, next) { //
10        getRawBody(req, { //
11            expected: req.headers['content-length'], //
12            limit: limit
13        }, >> anonymous_1000);
14    };
15 }
16
17 app.post('/image/text', bodyParser(1 * 1024 * 1024), routes.writeTextToImages); //
18 app.listen(8000);
19
20 flx anonymous_1000
21 -> express_dispatcher
22 function (err, buffer) { //
23     req.body = buffer; //
24     next_placeholder(req, -> express_dispatcher); //
25 }
26
27 flx express_dispatcher & express // 
28 -> null
29 merge(req, msg.req);
30 next(); //

```

Listing 5.5 – Simplified modification on the compiled result

Originally, the function `next` is the continuation to allow the anonymous callback on line 11, to continue the execution with the `next` function to handle the request. To isolate the anonymous callback, this function is replaced on both ends. The result of this replacement is illustrated in listing 5.5. The `express Router` registers a fluxion named `express_dispatcher`, line 27, to continue the execution after the fluxion `anonymous_1000`. This fluxion is in the same group `express` as the `main` fluxion, hence it has access to network sockets, to the original variable `req`, and to the original function `next`. The call to the original `next` function in the anonymous callback is replaced by a placeholder to push the stream to the fluxion `express_dispatcher`, line 24. The fluxion `express_dispatcher` receives the stream from the upstream fluxion `anonymous_1000`, merges back the modification in the variable `req` to propagate the side effects, before calling the original function `next` to continue the execution, line 30.

After the modifications detailed above, the server works as expected for the subset of functionalities we modified. The isolated fluxion correctly receives, and returns its serialized messages. The client successfully receives a gif frame containing the text.

5.3.3 Future works

We intend to implement the compilation process presented into the runtime. A just-in-time compiler would allow to identify callbacks dynamically evaluated, and to analyze the memory to identify side-effects propagations instead of relying only on the source code. Moreover, this memory analysis would allow the closure serialization required to compile application using higher-order functions.

Chapter 6

Conclusion

Appendix A

Language popularity

A.1 PopularitY of Programming Languages (PYPL)

¹ The PYPL index uses Google trends² as a leading indicator of the popularity of a programming language. It search for the trend for each programming language by counting the number of searches of this language and the word "tutorial".

PYPL for May 2015

¹<http://pypl.github.io/PYPL.html>

²<https://www.google.com/trends/>

Rank	Change	Language	Share	Trend
1		Java	24.1%	-0.9%
2		PHP	11.4%	-1.6%
3		Python	10.9%	+1.3%
4		C#	8.9%	-0.7%
5		C++	8.0%	-0.2%
6		C	7.6%	+0.2%
7		Javascript	7.1%	-0.6%
8		Objective-C	5.7%	-0.2%
9		Matlab	3.1%	+0.1%
10	2×↑	R	2.8%	+0.7%
11	5×↑	Swift	2.6%	+2.9%
12	1×↓	Ruby	2.5%	+0.0%
13	3×↓	Visual Basic	2.2%	-0.6%
14	1×↓	VBA	1.5%	-0.1%
15	1×↓	Perl	1.2%	-0.3%
16	1×↓	lua	0.5%	-0.1%

A.2 TIOBE

³

The TIOBE index uses many search engines as an indicator of the current popularity of programming languages. It counts the number of pages each search engine finds when queried with the language name and the word "programming". This indicator indicates the number of resources available, and the discussions about a given programming language.

Javascript was the most rising language of 2014 in the TIOBE index.
TIOBE for April 2015

³<http://www.tiobe.com/index.php/content/paperinfo/tpci/index.html>

Apr 2015	Apr 2014	Change	Programming Language	Ratings	Change
1	2	↑	Java	16.041%	-1.31%
2	1	↓	C	15.745%	-1.89%
3	4	↑	C++	6.962%	+0.83%
4	3	↓	Objective-C	5.890%	-6.99%
5	5		C#	4.947%	+0.13%
6	9	↑	JavaScript	3.297%	+1.55%
7	7		PHP	3.009%	+0.24%
8	8		Python	2.690%	+0.70%
9	-	2×↑	Visual Basic	2.199%	+2.20%

A.3 Programming Language Popularity Chart

⁴

The programming language popularity chart indicates the activity of a given language in the online communities. It uses two indicators to rank languages : the number of line changed in github of, and the number of questions tagged with a certain language.

Javascript is ranked number one in this index. The Javascript community is particularly active online, and in the open source.

indeed.com

A.4 Black Duck Knowledge

⁵

The black-duck, which analyze the usage of language on many forges, and collaborative hosts, rank Javascript number 2, after C, and with about the same usage as C++.

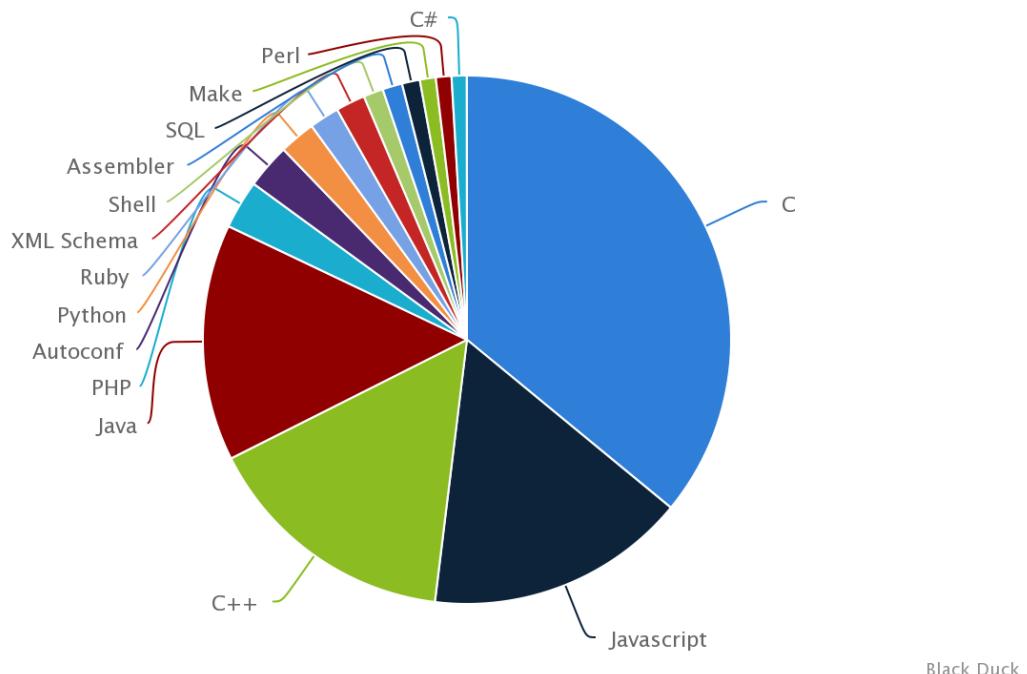
github.com sourceforge.net cpan.org rubyforge7.org planetsourcecode.com ddj.com

⁴<http://langpop.corger.nl>

⁵<https://www.blackducksoftware.com/resources/data/this-years-language-use>

Language	%
C	34.80
Javascript	15.45
C++	15.13
Java	14.02
PHP	2.87
Autoconf	2.65
Python	2.15
Ruby	1.77
XML Schema	1.73
Shell	1.18
Assembler	1.16
SQL	1.07
Make	0.94
Perl	0.92
C#	0.90

Releases within the last 12 months



A.5 Github

<http://githut.info/>

A.6 HackerNews Poll

<https://news.ycombinator.com/item?id=3746692>

Language	Count
Python	3335
Ruby	1852
JavaScript	1530
C	1064
C#	907
PHP	719
Java	603
C++	587
Haskell	575
Clojure	480
CoffeeScript	381
Lisp	348
Objective C	341
Perl	341
Scala	255
Scheme	202
Other	195
Erlang	171
Lua	150
Smalltalk	130
Assembly	116
SQL	112
Actionscript	109
OCaml	88
Groovy	83
D	79
Shell	76
ColdFusion	51
Visual Basic	47
Delphi	45
Forth	41
Tcl	34
Ada	29
Pascal	28
Fortran	26
Rexx	13
Cobol	12

Bibliography

- [1] D Abadi, D Carney, and U Cetintemel. “Aurora: a data stream management system”. In: *Proceedings of the ...* (2003).
- [2] DJ Abadi, Y Ahmad, and M Balazinska. “The Design of the Borealis Stream Processing Engine.” In: *CIDR* (2005).
- [3] DJ Abadi and D Carney. “Aurora: a new model and architecture for data stream management”. In: *The VLDB Journal— ...* (2003).
- [4] A Adya, J Howell, and M Theimer. “Cooperative Task Management Without Manual Stack Management.” In: *USENIX Annual Technical Conference* (2002).
- [5] T Akidau and A Balikov. “MillWheel: Fault-Tolerant Stream Processing at Internet Scale”. In: *Proceedings of the VLDB Endowment 6.11* (2013).
- [6] SP Amarasinghe, JAM Anderson, MS Lam, and CW Tseng. “An Overview of the SUIF Compiler for Scalable Parallel Machines.” In: *PPSC* (1995).
- [7] GM Amdahl. “Validity of the single processor approach to achieving large scale computing capabilities”. In: *Proceedings of the April 18-20, 1967, spring joint ...* (1967).
- [8] LO Andersen. “Program analysis and specialization for the C programming language”. In: (1994).
- [9] H Balakrishnan and M Balazinska. “Retrospective on aurora”. In: *The VLDB Journal* (2004).
- [10] U Banerjee. *Loop parallelization*. 2013.
- [11] JR von Behren, J Condit, and EA Brewer. “Why Events Are a Bad Idea (for High-Concurrency Servers).” In: *HotOS* (2003).
- [12] R Von Behren, J Condit, and F Zhou. “Capriccio: scalable threads for internet services”. In: *ACM SIGOPS ...* (2003).

- [13] M Bodin and A Chaguéraud. “A trusted mechanised JavaScript specification”. In: *Proceedings of the 41st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (2014).
- [14] I Buck, T Foley, and D Horn. “Brook for GPUs: stream computing on graphics hardware”. In: *... on Graphics (TOG)* (2004).
- [15] S Chandrasekaran and O Cooper. “TelegraphCQ: continuous dataflow processing”. In: *Proceedings of the ...* (2003).
- [16] J Chen, DJ DeWitt, F Tian, and Y Wang. “NiagaraCQ: A scalable continuous query system for internet databases”. In: *ACM SIGMOD Record* (2000).
- [17] Austin T. Clements, M. Frans Kaashoek, Nickolai Zeldovich, Robert T. Morris, and Eddie Kohler. “The scalable commutativity rule”. In: *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles - SOSP '13*. New York, New York, USA: ACM Press, Nov. 2013, pp. 1–17. DOI: [10.1145/2517349.2522712](https://doi.org/10.1145/2517349.2522712).
- [18] C Consel, H Hamdi, and L Réveillère. “Spidle: a DSL approach to specifying streaming applications”. In: *Generative ...* (2003).
- [19] Melvin E. Conway. “Design of a separable transition-diagram compiler”. In: *Communications of the ACM* 6.7 (July 1963), pp. 396–408. DOI: [10.1145/366663.366704](https://doi.org/10.1145/366663.366704).
- [20] J Dean and S Ghemawat. “MapReduce: simplified data processing on large clusters”. In: *Communications of the ACM* (2008).
- [21] Edsger W. Dijkstra. “Guarded commands, nondeterminacy and formal derivation of programs”. In: *Communications of the ACM* 18.8 (Aug. 1975), pp. 453–457. DOI: [10.1145/360933.360975](https://doi.org/10.1145/360933.360975).
- [22] Edsger W. Dijkstra. “Letters to the editor: go to statement considered harmful”. In: *Communications of the ACM* 11.3 (Mar. 1968), pp. 147–148. DOI: [10.1145/362929.362947](https://doi.org/10.1145/362929.362947).
- [23] Edsger W. Dijkstra. “The structure of the “THE”-multiprogramming system”. In: *Communications of the ACM* 11.5 (May 1968), pp. 341–346. DOI: [10.1145/363095.363143](https://doi.org/10.1145/363095.363143).
- [24] Raul Castro Fernandez, Matteo Migliavacca, Evangelia Kalyvianaki, and Peter Pietzuch. “Making state explicit for imperative big data processing”. In: *USENIX ATC* (2014).
- [25] JI Fernández-Villamor. “Microservices-Lightweight Service Descriptions for REST Architectural Style.” In: *... 2010-Proceedings of ...* (2010).

- [26] D Flanagan. *JavaScript: the definitive guide*. 2006.
- [27] P Gardner and G Smith. “JuS: Squeezing the sense out of javascript programs”. In: *JSTools@ ECOOP* (2013).
- [28] PA Gardner, S Maffeis, and GD Smith. “Towards a program logic for JavaScript”. In: *ACM SIGPLAN Notices* (2012).
- [29] JJ Garrett. “Ajax: A new approach to web applications”. In: (2005).
- [30] SD Gribble, M Welsh, and R Von Behren. “The Ninja architecture for robust Internet-scale systems and services”. In: *Computer Networks* (2001).
- [31] NJ Gunther. “A New Interpretation of Amdahl’s Law and Geometric Scalability”. In: *arXiv preprint cs/0210017* (2002).
- [32] NJ Gunther. “A simple capacity model of massively parallel transaction systems”. In: *CMG-CONFERENCE-* (1993).
- [33] NJ Gunther. “Understanding the MP effect: Multiprocessing in pictures”. In: *In other words* (1996).
- [34] JL Gustafson. “Reevaluating Amdahl’s law”. In: *Communications of the ACM* (1988).
- [35] CT Haynes, DP Friedman, and M Wand. “Continuations and coroutines”. In: *... of the 1984 ACM Symposium on ...* (1984).
- [36] B He, M Yang, Z Guo, R Chen, and B Su. “Comet: batched stream processing for data intensive distributed computing”. In: *... on Cloud computing* (2010).
- [37] C Hewitt. “Viewing control structures as patterns of passing messages”. In: *Artificial intelligence* (1977).
- [38] C Hewitt, P Bishop, and R Steiger. “A universal modular actor formalism for artificial intelligence”. In: *Proceedings of the 3rd international joint conference on Artificial intelligence* (1973).
- [39] C. A. R. Hoare. “Communicating sequential processes”. In: *Communications of the ACM* 21.8 (Aug. 1978), pp. 666–677. DOI: [10.1145/359576.359585](https://doi.org/10.1145/359576.359585).
- [40] C. A. R. Hoare. “Monitors: an operating system structuring concept”. In: *Communications of the ACM* 17.10 (Oct. 1974), pp. 549–557. DOI: [10.1145/355620.361161](https://doi.org/10.1145/355620.361161).
- [41] YW Huang, F Yu, C Hang, and CH Tsai. “Securing web application code by static analysis and runtime protection”. In: *Proceedings of the 13th ...* (2004).

- [42] M Isard, M Budiu, Y Yu, A Birrell, and D Fetterly. “Dryad: distributed data-parallel programs from sequential building blocks”. In: *ACM SIGOPS Operating ...* (2007).
- [43] D Jang and KM Choe. “Points-to analysis for JavaScript”. In: *Proceedings of the 2009 ACM symposium on Applied ...* (2009).
- [44] N Jovanovic, C Kruegel, and E Kirda. “Pixy: A static analysis tool for detecting web application vulnerabilities”. In: *Security and Privacy, 2006 ...* (2006).
- [45] Gilles Kahn. “The semantics of a simple language for parallel programming”. In: (1974).
- [46] Gilles Kahn and David Macqueen. *Couroutines and Networks of Parallel Processes*. en. Tech. rep. 1976, p. 20.
- [47] S Krishnamurthy and S Chandrasekaran. “TelegraphCQ: An architectural status report”. In: *IEEE Data Eng. ...* (2003).
- [48] MN Krohn, E Kohler, and MF Kaashoek. “Events Can Make Sense.” In: *USENIX Annual Technical Conference* (2007).
- [49] L Lamport. “Time, clocks, and the ordering of events in a distributed system”. In: *Communications of the ACM* (1978).
- [50] D Logothetis, C Olston, and B Reed. “Stateful bulk processing for incremental analytics”. In: *Proceedings of the 1st ...* (2010).
- [51] S Maffeis, JC Mitchell, and A Taly. “An operational semantics for JavaScript”. In: *Programming languages and systems* (2008).
- [52] S Maffeis, JC Mitchell, and A Taly. “Isolating JavaScript with filters, rewriting, and wrappers”. In: *Computer Security—ESORICS 2009* (2009).
- [53] WR Mark and RS Glanville. “Cg: A system for programming graphics hardware in a C-like language”. In: ... *Transactions on Graphics* (... (2003).
- [54] ND Matsakis. “Parallel closures: a new twist on an old idea”. In: *Proceedings of the 4th USENIX conference on Hot ...* (2012).
- [55] MD McCool. “Structured parallel programming with deterministic patterns”. In: *Proceedings of the 2nd USENIX conference on Hot ...* (2010).
- [56] M Migliavacca and D Eyers. “SEEP: scalable and elastic event processing”. In: *Middleware'10 Posters ...* (2010).
- [57] G Moore. “Cramming More Components Onto Integrated Circuits”. In: *Electronics* 38 (1965), p. 8.

- [58] JP Morrison. *Flow-Based Programming*. 1994, pp. 1–377.
- [59] JF Naughton, DJ DeWitt, and D Maier. “The Niagara internet query system”. In: *IEEE Data Eng.* ... (2001).
- [60] R Nelson. “Including queueing effects in Amdahl’s law”. In: *Communications of the ACM* (1996).
- [61] L Neumeyer and B Robbins. “S4: Distributed stream computing platform”. In: *Data Mining Workshops* ... (2010).
- [62] VS Pai, P Druschel, and W Zwaenepoel. “Flash: An efficient and portable Web server.” In: *USENIX Annual Technical Conference* (1999).
- [63] DL Parnas. “On the criteria to be used in decomposing systems into modules”. In: *Communications of the ACM* (1972).
- [64] R Power and J Li. “Piccolo: Building Fast, Distributed Programs with Partitioned Tables.” In: *OSDI* (2010).
- [65] Z Qian, Y He, C Su, Z Wu, and H Zhu. “Timestream: Reliable stream computation in the cloud”. In: *Proceedings of the 8th ACM European Conference on Computer Systems (EuroSys ’13)* (2013).
- [66] C Radoi, SJ Fink, R Rabbah, and M Sridharan. “Translating imperative code to MapReduce”. In: *Proceedings of the 2014* ... (2014).
- [67] DP Reed. “” Simultaneous” Considered Harmful: Modular Parallelism.” In: *HotPar* (2012).
- [68] MC Rinard and PC Diniz. “Commutativity analysis: A new analysis framework for parallelizing compilers”. In: *ACM SIGPLAN Notices* (1996).
- [69] Tiago Salmito, Ana Lucia de Moura, and Noemi Rodriguez. “A Flexible Approach to Staged Events”. English. In: *2013 42nd International Conference on Parallel Processing*. IEEE, Oct. 2013, pp. 661–670. DOI: [10.1109/ICPP.2013.80](https://doi.org/10.1109/ICPP.2013.80).
- [70] Tiago Salmito, Ana Lúcia de Moura, and Noemi Rodriguez. “A stepwise approach to developing staged applications”. In: *The Journal of Supercomputing* (Jan. 2014). DOI: [10.1007/s11227-014-1110-4](https://doi.org/10.1007/s11227-014-1110-4).
- [71] GD Smith. “Local reasoning about web programs”. In: (2011).
- [72] M Sridharan, J Dolby, and S Chandra. “Correlation tracking for points-to analysis of JavaScript”. In: *ECOOP 2012—Object-* ... (2012).

- [73] W. P. Stevens, G. J. Myers, and L. L. Constantine. “Structured design”. English. In: *IBM Systems Journal* 13.2 (1974), pp. 115–139. DOI: [10.1147/sj.132.0115](https://doi.org/10.1147/sj.132.0115).
- [74] GJ Sussman and GL Steele Jr. “Scheme: A interpreter for extended lambda calculus”. In: *Higher-Order and Symbolic Computation* (1998).
- [75] W Thies, M Karczmarek, and S Amarasinghe. “StreamIt: A language for streaming applications”. In: *Compiler Construction* (2002).
- [76] A Toshniwal and S Taneja. “Storm@ twitter”. In: *Proceedings of the 2014 ACM SIGMOD international conference on Management of data - SIGMOD ’14* (2014).
- [77] S Wei and BG Ryder. “State-sensitive points-to analysis for the dynamic behavior of JavaScript objects”. In: *ECOOP 2014—Object-Oriented Programming* (2014).
- [78] M Welsh, D Culler, and E Brewer. “SEDA: an architecture for well-conditioned, scalable internet services”. In: *ACM SIGOPS Operating Systems Review* (2001).
- [79] Sunny Wong, Yuanfang Cai, Giuseppe Valetto, Georgi Simeonov, and Kanwarpreet Sethi. “Design Rule Hierarchies and Parallelism in Software Development Tasks”. In: *2009 IEEE/ACM International Conference on Automated Software Engineering*. IEEE, Nov. 2009, pp. 197–208. DOI: [10.1109/ASE.2009.53](https://doi.org/10.1109/ASE.2009.53).
- [80] D Yu, A Chander, N Islam, and I Serikov. “JavaScript instrumentation for browser security”. In: *ACM SIGPLAN Notices* (2007).
- [81] Yuan Yu, Michael Isard, Dennis Fetterly, Mihai Budiu, Ulfar Erlingsson, Pradeep Kumar Gunda, Jon Currey, Frank McSherry, Kannan Achan, and Christophe Poulain. “Some sample programs written in DryadLINQ”. In: *Microsoft Research* (2009).