UNIVERSITÉ DE MONTRÉAL

TITRE DE MON DOCUMENT

JULES CÆSAR
DÉPARTEMENT DE GÉNIE INFORMATIQUE ET GÉNIE LOGICIEL
ÉCOLE POLYTECHNIQUE DE MONTRÉAL

THÈSE PRÉSENTÉE EN VUE DE L'OBTENTION
DU DIPLÔME DE PHILOSOPHIÆ DOCTOR
(GÉNIE INFORMATIQUE)
AVRIL 2010

UNIVERSITÉ DE MONTRÉAL


ÉCOLE POLYTECHNIQUE DE MONTRÉAL



Cette thèse intitulée :


TITRE DE MON DOCUMENT




présentée par : CÆSAR Jules
en vue de l'obtention du diplôme de : Philosophiæ Doctor
a été dûment acceptée par le jury d'examen constitué de :




M. NOM Prénom, Doct., président
Mme NOM Prénom, Ph. D., membre et directrice de recherche
M. NOM Prénom, Ph. D., membre

# DEDICATION

*To all friends at the lab,*
*I will miss you. . .*

# ACKNOWLEDGEMENTS

In the acknowledgements, the author points out the help that various people have provided, including advice or any other type of contribution as the author carried out their research. As appropriate, this is the section where the candidate must thank their thesis or dissertation supervisor, grant-awarding organizations, or companies that provided bursaries or research funds.

If the thesis or dissertation is in French and students wish to thank someone in particular in English, they must insert a second page and title it in English (i.e., "Acknowledgments"). If the thesis or dissertation is in English and students wish to thank someone in particular in French, they must insert a second page and title it in French (i.e., "Remerciements").

# RÉSUMÉ

The résumé (mandatory summary) is a brief explanation in French of the work's topic, its objectives, the research questions or hypotheses put forth, the experimental methods used, and the results analysis. It also includes the key research conclusions and future applications. In general, a summary does not exceed three pages.

The summary must provide an exact idea of the thesis or dissertation's content. It cannot be a simple enumeration of the manuscript's parts. The goal is to precisely and concisely present the nature and scope of the research. A summary must never include references or figures.

# ABSTRACT

Written in English, the abstract is a brief summary similar to the previous section (Résumé). However, this section is not a word for word translation of the French.

# TABLE OF CONTENTS

# LIST OF TABLES

# LIST OF FIGURES

# LIST OF SYMBOLS AND ABBREVIATIONS

IETF   Internet Engineering Task Force

OSI   Open Systems Interconnection

# LIST OF APPENDIXES

# CHAPTER 1    INTRODUCTION

Texte en *italique*, PETITES MAJUSCULES, mot insécable.

Texte souligné, <mark>surligné</mark>, **gras**.

Texte entre "guillemets".

Police `monospace`.

Un mot courant en réseautique mobile : nœud[1].

L'objet RSVP `SENDER_TEMPLATE`.

Nom d'un auteur : Postel.

Une architecture 32 bits.

## 1.1    Définitions et concepts de base

1$^{\text{re}}$ utilisation d'un acronyme : Internet Engineering Task Force (IETF).

2$^{\text{e}}$ utilisation d'un acronyme : IETF.

Acronyme au long : Internet Engineering Task Force.

### 1.1.1    Une sous-section

Un URL : École Polytechnique de Montréal.

#### Une sous-sous-section

Les besoins des flots de données peuvent être catégorisés selon quatre paramètres importants (voir Tanenbaum, 2002, sect. 5.4) ou :

— la fiabilité (acheminement des données avec succès) ;

— le délai de bout-en-bout de la source vers la destination ;

— la variation du délai de bout-en-bout (*jitter*) ;

— la bande passante requise (le débit des informations).

**Le niveau paragraphe**   est plus bas encore dans la hiérarchie... Une citation entre parenthèses (voir Zhang et al., 2005). ou des citations entre parenthèses (Nichols et al., 2010; Abondo, 2005; Tseng et al., 2003).

---

1. Note de bas de page.

## 1.2 Éléments de la problématique

La description de l'en-tête commun de RSVP est détaillée ci-dessous :

Ver :      4 bits
           Version du protocole. La version actuelle est 1.

Flags :    4 bits
           Aucun Flag n'est défini. L'émetteur doit (**MUST**) mettre le
           champ à zéro et le récepteur doit (**MUST**) ignorer ce champ.

Msg Type : 8 bits
           Type de message

Checksum : 16 bits
           Complément à un du complément à un de la somme des champs
           de l'en-tête, avec le champ Checksum à 0 pour des fins de calcul.
           La valeur 0 signifie qu'aucun Checksum n'a été transmis. Si le
           résultat du calcul du Checksum donne 0, la valeur 0xFFFF doit
           être stockée dans ce champ.

TTL :      8 bits
           Valeur originelle du champ TTL utilisée pour transmettre ce
           message.

Reserved : 8 bits
           Réservé pour usage futur. L'émetteur doit (**MUST**) mettre le
           champ à zéro et le récepteur doit (**MUST**) ignorer ce champ.

Length :   16 bits
           Longueur totale du message en octets, incluant l'en-tête commun
           et tous les objets de longueur variable.

### 1.2.1 Autres types de structures de données

L'énumération :

1. Un item ;
2. Un autre item.

### 1.2.2 Le protocole IPv6

Voir la Figure 1.1 pour plus de détails. Le champs DSCP est décrit dans le Tableau 1.1.

Figure 1.1 L'en-tête IPv6

Table 1.1 Plages de valeurs pour le champ `DSCP`

| Plage | Valeurs | Règle d'assignation |
|-------|---------|---------------------|
| 1 | xxxxx0 | Assignation par une norme de l'IANA |
| 2 | xxxx11 | Expérimentation/Usage local |
| 3 | xxxx01 | Expérimentation/Usage local (pourrait être jointe à la plage 1) |

## 1.3 Objectifs de recherche

Les objectifs de la recherche sont de concevoir un algorithme $O(n)$.

## 1.4 Plan du mémoire

Un tableau :

La formule d'Erlang-B :

$$P_b = \frac{\frac{\rho^C}{C!}}{\sum\limits_{x=0}^{C} \frac{\rho^x}{x!}} \tag{1.1}$$

Une autre équation :

$$\begin{aligned}
P_c &= (1 - P_b) \times (1 - P_f)^N \\
&= (1 - P_b)^{N+1}
\end{aligned} \tag{1.2}$$

Table 1.2 Constantes et variables du modèle analytique

| Symbole | Description |
|---|---|
| $\lambda$ | Taux d'arrivée moyen des requêtes de réservation de ressources |
| $\frac{1}{\mu}$ | Durée moyenne d'une session |
| $C$ | Capacité d'une cellule (nombre de sessions supportées) |
| $v_{moy}$ | Vitesse moyenne des MN dans le réseau d'accès |
| $L$ | Longueur d'un côté d'une cellule carrée |
| $n$ | Nombre moyen de MN dans une cellule |
| $\rho$ | Charge d'une cellule |
| $P_b$ | Probabilité de blocage d'une requête de réservation |
| $P_f$ | Probabilité d'interruption forcée d'une session |
| $P_c$ | Probabilité de compléter une session avec succès |
| $\Delta T$ | Délai de transmission |

Enfin, l'expression suivante indique le moment à partir duquel les réservations de ressources sont en place :

$$\Delta T_{init} = \begin{cases} 2\Delta T_{E2E} & \Delta T_{wan} > (\Delta T_{rad} + \Delta T_{net}) \\ \Delta T_{E2E} + 3(\Delta T_{rad} + \Delta T_{net}) & \text{sinon} \end{cases} \tag{1.3}$$

**Le taux de paquets perdus** correspond au nombre de paquets éliminés à cause d'une erreur de *checksum* à un nœud quelconque ou d'une situation de congestion. Le taux de paquets perdus pour un chemin est déterminé de la façon suivante :

$$PLR_P = 1 - \prod_{i=1}^{N}(1 - PLR_i) \tag{1.4}$$

Toutefois, si les taux d'erreurs sont très faibles, comme c'est généralement le cas pour des liens optiques, on peut approximer $PLR_P$ de façon à le transformer en un paramètre additif :

$$\begin{aligned} PLR_{L_1 \oplus L_2} &= 1 - (1 - PLR_1)(1 - PLR_2) \\ &= 1 - (1 - PLR_2 - PLR_1 + \underbrace{PLR_1 \times PLR_2}_{\text{négligeable}}) \qquad PLR_1 \ll 1, PLR_2 \ll 1 \\ &\approx PLR_1 + PLR_2 \end{aligned} \tag{1.5}$$

Une courbe :



Figure 1.2 Délai moyen en fonction du taux d'utilisation d'un lien

This paragraph is formatted by LaTeX according to the standard rules of the English language (e.g. hyphenation).

L'arithmétique en virgule flottante peut entraîner des erreurs d'approximation et il est important d'en être conscient (voir Goldberg, 1991).

De même, les calculs effectués sur une carte graphique (GPU) peuvent introduire des erreurs d'approximation (Benz et al., 2012; D'Silva et al., 2012; Dabrowski et al., 2011; De Dinechin et al., 2011; de Figueiredo and Stolfi, 2004; Filliâtre and Marché, 2007; Fousse et al., 2007; Goubault, 2001; Goubault et al., 2008; Harder and Khoury, 2005; Higham, 2002; Tanenbaum, 2002; Whitehead and Fit-Florea, 2011; Johansson et al., 2013; Nichols et al., 2010; NVIDIA, 2012; Benz et al., 2012; Bao and Zhang, 2013).

# CHAPTER 2     CRITICAL LITERATURE REVIEW

Texte.

**CHAPTER 3    From Academic Research to Open Source Software**

Text.

# CHAPTER 4    SECOND TOPIC

As software evolves, a developer's contributions will gradually vanish as they are being replaced by other developers' code, slowly eroding the developer's footprint in the software project. Even though this developer's knowledge of the file did not disappear overnight, to outsiders, the developer and her expertise have become invisible. Through an empirical study on 5 years of Linux development history, this paper analyses this phenomenon of expertise erosion by building a 2-dimensional model of developer expertise involving a range of developer activities and involving activity data on more than one release. Using these models, we found that although many Linux maintainers' own coding footprint has regressed over time, their expertise is perpetuated through involvement in other development activities such as patch reviews and committing upstream on behalf of other developers. Considering such activities over time further improves the expertise models.

## 4.1    Introduction

As reported by Damien et al. **?**, employee turnover is a major challenge of information technology organizations. Estimations of the cost of losing an employee amount to between 0.5 and 1.5 times her salary, with the cost of replacing a software engineer in particular exceeding $100,000 **?**. These costs are not limited to the software engineer's company, but also spread to open source development. In their 2017 Linux kernel report, Corbet et al. **?** noted that "well over 85 percent of all kernel development is demonstrably done by developers who are being paid for their work". In fact, only 8.2% of all kernel commits were made by volunteer developers. Hence, developer turn-over in companies risks to impact open source development as well!

Apart from improving the working conditions and onboarding procedures, software organizations (both closed and open source) need to invest time in finding the "right" expert to replace a parting developer. While it is possible to train newcomers and bring them up to speed (e.g., one third of the kernel contributors in the last 12 months were newcomers, and 60% of those made their contribution as employee), the term "right" refers to having a similar profile, allowing the new expert to seamlessly fit in and continue his or her predecessor's work, without significant loss of knowledge. Thanks to the widespread adoption of agile development and open source development models, software development has become a collaborative endeavor, in which knowledge is shared across the members of an organization, hence in principle it should be possible to find contributors with similar profiles.

Unfortunately, there is no consensus on how to measure the profile of a developer, and how to determine whether such a profile indicates the developer to be an expert. The simplest way to measure someone's development activities is to count the number of code changes (e.g., Git commits) authored. This is for example how Corbet et al. determine the most active developers and organizations in Linux kernel development **?**. Yet, at the same time, they note that "The total number of patches signed off by Linus Torvalds (207, or 0.3 percent of the total) continues its long-term decline. That reflects the increasing amount of delegation to subsystem maintainers who do the bulk of the patch review and merging." Worse, developer rankings based on the number of commits differed substantially based on the period during which this metric was measured (e.g., ranking based on last 10 years vs. last year). At a minimum, one needs to be careful interpreting these and other measures such as a developer's code legacy as shown by "git blame" **????**.

To make developer expertise measures more robust and reliable, this paper proposes a 2-dimensional developer expertise footprint, addressing two important issues with current expertise models. First of all, while the amount of code written by a person can be an important indicator of expertise, it does not take into consideration the actions of people who do not directly contribute source code, such as those who review it or discuss it on mailing lists or issue repositories. Our footprint measure combines indicators of multiple kinds of developer activities.

Second, as indicated by the Corbet et al., current measures focus on a given software release or development period, basically ignoring the development activities that happened before. While a person who wrote 50% of the code changes of the previous release could be less of an expert than a person who wrote 50% of the code changes of the current release, the former developer might have been ill or absent, or might have been the one mentoring the latter developer. As such, both developers should be considered as experts, not just the latter developer. Our footprint measure allows to consider a developer's activities across different time periods.

We empirically evaluate the expertise footprint models on 5 years of Linux kernel development history, addressing the following research questions :

***RQ1)*** *How does expertise evolve in an open source project?*

> Almost 1 out of 4 subsystems has seen a change in maintainership during the last 22 Linux kernel releases, with the code footprint of maintainers gradually decreasing over time.

***RQ2)*** *How well does dimension 1 explain expertise?*

Models involving a maintainers' own code footprint and coordination activities (committing and/or reviewing) perform the best.

***RQ3)*** *How well does dimension 2 explain expertise?*

Expertise models considering the last $R$ releases perform better than single-release models.

## 4.2  Background and Related Work

Maintainers ensure the longevity of the Linux kernel by not only contributing new code, but also by reviewing and integrating code submitted to their subsystems by other developers. They are the backbone of Linux kernel development, yet few studies have been done to study their work **?**.

In particular, a maintainer's departure of her subsystem calls for her immediate replacement. However, only a developer with extensive experience in the subsystem can take on the task of maintainer. Software expertise and knowledge have been extensively studied in the past **????**. Many different models were created to attempt to assess developer expertise.

Earlier expertise models **??** made the assumption that expertise is related to coding activities. They both extracted data describing the changes in the source code from the version control system. In other words, they measured a developer's expertise in terms of the number of changes made to the system.

Fritz et al. **?** later expressed their concern about the assumption that activity indicates expertise. Their review of psychology studies indicated a lack of evidence proving that activity can be an indicator of knowledge. However, after a qualitative study consisting of 19 java developers interviews, they were able to confirm a relationship between commit frequency and expertise. In addition to that, they foud evidence proving that authorship (as obtained from the amount of churn contributed, or through "git blame") is also capable of indicating expertise.

Bhattacharya et al. **?** explored the suitability of different expertise indicators depending on the developer's role. They argue that state-of-the-art metrics (lines of code and commits added), being unaware of the developer's role, can lead to inaccurate results. They add that code activity metrics like the number of lines of code added, only describe expertise at a local level and poorly capture global expertise.

Even though the goal of each study mentioned above varies, they all introduced expertise detection models relying on two simple metrics to measure expertise : the number of lines of

code contributed and/or the number of commits.

Although these state-of-the-art techniques are well suited to detect experts among regular developers, we believe that they do not properly evaluate more complex expertise, such as that of subsystem maintainers. Most of the daily tasks, such as reviewing, of maintainers are ignored, creating an inherent bias in the expertise models.

We address this bias by building expertise models based on a variety of metrics capturing the full breadth of software development activities, and also considering a longer the evolution of such activities over time.

## 4.3 Measuring the Expertise Footprint of Contributors

This section discusses the two-dimensional model of expertise footprint proposed by this paper to enable identification of experienced team members (e.g., developers, testers, etc.). The first dimension of the footprint model considers a wide range of activities performed by a project member, not only focusing on code changes, but also code review or even a developer's code "legacy" (i.e., contributed code that still survives in the code base). The second dimension enhances the first dimension by not only considering the range of activities in the latest release, but *across the last N releases*. As such, accidental lulls or shifts in project activity are accounted for.

Note that the expertise we are interested in is expertise about the *internals* of a particular source code file or component. An alternative form of expertise would consider knowledge on how to *use* a particular component (API). We focus on the former kind of expertise, since it is at the heart of a software organizations needs. For example, it allows to measure the expertise of a particular individual, allowing the organization to better use her skills, evaluate her value to the organization, and assess the risk of her potential departure. Furthermore, it is important for an organization to know—for any section of the system—who are its experts, and their level of expertise. Finally, in both cases, it is also important to know how this expertise is changing over time (e.g., the areas where a person is gaining and losing expertise, and, for a given area, how it is losing or gaining experts).

### 4.3.1 Dimension 1 : Contributor Activities

The expertise footprint model that we propose explicitly considers a wide range of development activities instead of focusing only on review- or code-related activities. Table 4.1 provides a non-exhaustive list of activities, from very technical to outreach activities. Any activity by a contributor to one of these, can increase (or at least maintain) the contributor's

| activity | definition |
|---|---|
| `legacy` | influential code contributed by C that still survives in $R_i$ |
| `authored` | code authored by C since $R_{i-1}$ |
| `committed` | code committed by C since $R_{i-1}$ |
| `reviewed` | code changes reviewed by C since $R_{i-1}$ |
| `translated` | involvement in translating/localizing textual strings for $R_i$ |
| `integrated` | effort spent by C integrating code changes since $R_{i-1}$ |
| `discussed` | effort spent by C discussing issue reports since $R_{i-1}$ |
| `represented` | effort spent by C representing S on social media since $R_{i-1}$ |
| `planned` | effort spent by C planning $R_i$ |

Table 4.1 Non-exhaustive list of activity measures that can be measured for a particular contributor C of a specific subsystem S in a given release $R_i$.

knowledge about the subsystem she is working in. The more measures are considered, the more comprehensive the expertise footprint model becomes, hence the better the expected performance for identification of experts in a project under study, provided the activities are weighted based on their relevance for a given project.

This flexibility comes at the expense of additional effort for mining these activity measures. Fortunately, when developers contribute code to an open or closed source project, data about each code change, code review or other activity is automatically stored in the project's software repositories. The most trivial example are the code changes (commits) recorded in a version control system like git. However, information about the contributor's activity in issue report discussions is also readily available from the project's issue repository (e.g., bugzilla or jira), code review activity from the review repository (e.g., gerrit or mailing list) and mailing list activity from the mailing list archive. Of the metrics in Table 4.1, `represented` and `planned` are the hardest to obtain data for.

Given a set of activity measures $\mathbb{A} = \{a_i | a_i \ is \ activity \ measure\}$, we compute the expertise footprint of a release $j$ as :

$$footprint_j(\mathbb{A}) = \sum_i \frac{w_i \times a_i}{a_i^{tot}}$$

, where $w_i$ is a weight factor given to $a_i$ ($\sum_i w_i = 1$) and $a_i^{tot}$ is the total number of activities of a given type (e.g., number of source code lines, commits or reviews) recorded for a given activity and release. In other words, each activity is normalized, and the weighted sum of the normalized activities yields the $footprint_j(\mathbb{A})$ percentage. Hence, to instantiate the generic $footprint_j(\mathbb{A})$ measure, an organization first has to select the activity measures $\mathbb{A}$ relevant to its context, then determine the relative weight $w_i$ of each selected activity.

It is necessary to normalize each activity's measure to provide a better understanding of the true impact of developers' contributions in the subsystem. This is because the studied subsystems differ in size and the heuristic counts are inherently uneven by nature. For example, the value for `legacy` (in number of lines of code) will likely be much larger than the values of `authored` or `reviewed` (in number of commits).

### 4.3.2 Dimension 2 : Historical Perspective

While the definition of $footprint_j(\mathbb{A})$ takes into account a wide range of activities, it only considers a contributor's activity for one specific release $j$. As such, this measure might still provide misleading information when used to find the most appropriate expert for a given subsystem (e.g., to help debug a coding issue).

First of all, contributors in both closed and open source development evolve according to a particular career path. Even in open source, many contributors start out translating textual strings, before contributing smaller code fixes and ever larger changes until they are trusted enough to be able to review or even accept other contributors' code. This not only implies that a contributor's volume of contributions is scattered across different activities, but also that this scattering (and volume) might change over time. Hence, depending on the release under study, different $footprint_j(\mathbb{A})$ values are obtained, as if a specific contributor suddenly would have "lost" or "gained" a substantial percentage of expertise (footprint). To counter this noise, one should incorporate past experience to obtain a more robust footprint model.

Second, even when a contributor's responsibilities are stable across a time period, accidental life events such as illness or busy periods at work, or project events such as the scope of the upcoming release (major release vs. bug fix release) could lead to increases or decreases for certain activities. Again, if the contributor was an expert in the previous release, she will not have lost all of this expertise in one release due to illness. Hence, a release-specific $footprint_j(\mathbb{A})$ measure again would yield the wrong impression.

For this reason, the second dimension of our footprint model explicitly takes into account history by taking the weighted sum of $footprint_j(\mathbb{A})$ over the last R releases. In particular :

$$footprint_j^R(\mathbb{A}) = \sum_{i=j}^{j-R} W_i \times footprint_i(\mathbb{A})$$

, where $W_i$ is a weight factor given to the specific footprint of release $i$ ($\sum_i W_i = 1$). Note that $footprint_j(\mathbb{A}) = footprint_j^0(\mathbb{A})$, i.e., the footprint model obtained based on the first dimension is a special case of the second dimension ($R = 0$).

While the choice of weights $w_i$ for dimension 1 could be chosen arbitrarily based on relevance of individual activities, the weights $W_i$ typically will be decreasing, since recent activity typically is at least as important as older activity. For example, the weights could be linearly decaying (e.g., $[0.33, 0.27, 0.20, 0.13, 0.07]$ for $R = 4$), giving each older release proportionally less influence on the footprint model. Alternatively, an exponential ($[0.64, 0.23, 0.09, 0.03, 0.01]$) or logarithmic ($[0.34, 0.29, 0.23, 0.14, 0.00]$) decay could be used to give older release less or more influence, or (less likely) even a uniform ($[0.20, 0.20, 0.20, 0.20, 0.20]$) decay to give all considered releases the same importance.

### 4.3.3 Use Cases for Expertise Footprint

Given the footprint models $footprint_j(\mathbb{A})$ and $footprint_j^R(\mathbb{A})$, a number of use cases can be imagined.

The main use case considered in this paper is the identification of experts in a software project. Newcomers to a project typically are not aware who to ask for advice when encountering a specific technical issue. Similarly, when the maintainer of a specific component or library decides to retire, finding a good replacement is not always straightforward for an organization, as important development knowledge (across a range of development activities) risks to be lost.

A less straightforward application was suggested at one of the 2017 OPNFV Summit's panels, where substantial attention went to the issue of non-responsive Linux kernel maintainers. These are experts responsible for a given subsystem who, due to personal events, loss of interest or other reasons, start becoming non-responsive in communication with other developers or management. Having a reliable expertise measure in place would enable monitoring over time of maintainers' activities to spot long-term periods with sub-par performance. Such pro-active detection of issues could also suggest alternative maintainers.

Similarly, an expertise footprint can help an organization guard itself against accidental loss of manpower. For example, the bus factor **?** is a known measure of the risk that key personnel might disappear from a project, either out of free will or due to an accident. Organizations with a high bus factor could leverage an expertise footprint to identify backups for key developers or managers. As such, for each subsystem, an organization could have a list of the main people working in it as well as their expertise level.

In order to use the footprint models to find the most appropriate expert of a given subsystem, one needs to calculate $footprint_j(\mathbb{A})$ and/or $footprint_j^R(\mathbb{A})$ for each person who contributed at least once to one of the activities in $\mathbb{A}$. Then, the resulting footprint values should be ranked

from high to low. Ideally, the contributor with the highest footprint value is recommended as first candidate expert, followed by the contributor with the next highest footprint value, etc.

## 4.4 Case Study Setup

This section presents the design of an empirical study on the Linux kernel to evaluate the 2-dimensional footprint model introduced in the previous section. The study addresses the following research questions :

RQ1 : How does expertise evolve in an open source project?

RQ2 : How well does dimension 1 explain expertise?

RQ3 : How well does dimension 2 explain expertise?

### 4.4.1 Subject Data

Our study evaluates the expertise footprint models in the context of the Linux kernel. First of all, the Linux kernel is one of the hallmark open source projects, with a long history, large code base and vast supply of contributors. Second, the kernel is one of the few open source projects in which expertise is documented explicitly. The code base contains a file named `MAINTAINERS` that lists, for each subsystem, the experts in charge. Just as for source code, changes in maintainership are recorded through regular commits. This provides us with a unique oracle for our expertise measures.

Furthermore, the Linux Foundation (who governs and mentors the development of the Linux kernel and related open source initiatives) recently has started up the CHAOSS committee on Community Health Analytics for Open Source Software [1]. Amongst others, the aim of this committee is to identify explicit measures of expertise that can help prospective adopters of open source projects in choosing the right developer or maintainer to contact. As such, our study can help this concrete initiative, and we are in contact with the CHAOSS consortium.

Determining expertise footprints in the Linux kernel, especially taking into account the second dimension of our measure, requires a large set of historical data. We conducted our analysis on a set of 27 releases of the Linux kernel, spanning releases *v3.5* to *v4.11*, which corresponds to approximately 5 years of development and release history.

---

1. `https://chaoss.community/`

### 4.4.2 Filtering of the Data

Because of the constantly changing nature of the kernel, new subsystems are being added to the Linux kernel in every release to meet the demands associated with new hardware and changes in user expectations. Furthermore, it is not uncommon to see a subsystem disappearing, or, more precisely, becoming obsolete or orphaned (**?**, v4.11, `MAINTAINERS`, Line 84).

On the other hand, the importance of the historical aspect of our analysis forces us to choose long-standing subsystems that would best reflect the evolution of expertise of the subsystems' maintainers. Hence, we filter our Linux kernel data set to keep only subsystems that existed throughout the studied timespan. This subset reduces the number of subsystems from 1,662 to 734 subsystems.

For RQ2 and RQ3, we need further filtering to ensure a data set of subsystems for which there is ongoing activity in each studied release. To achieve this, we parsed, for each release, the `MAINTAINERS` file to extract each *active* subsystem along with its name, list of maintainer names, and the list of files and directories belonging to that subsystem.

We then retrieved the list of commits made to each subsystem, for each release that we considered. This allows us to compute, for each subsystem, the average number of commits across its releases. After matching each commit to its code reviews (see below), we also compute the average proportion of matched commits per release.

We then set minimum thresholds of 50 commits per release and 60% matched commits per release. This filtering reduces the 734 subsystems to a set of 78 subsystems for RQ2 and RQ3. This subset contains well know subsystems like `ARM PORT`, `XEN HYPERVISOR INTERFACE`, `SOUND`, `SCHEDULER`, and `CRYPTO API`.

### 4.4.3 Instantiation of the Footprint Models

Table 4.2 shows the five concrete activity measures considered in our empirical study on the Linux kernel. These measures cover influential source code contributed (`legacy`), the volume of code changes since the last release (`authored` and `committed`), and code review activities since the last release (`attributed` and `reviewed`).

Our $footprint_j(\mathbb{A})$ and $footprint_j^R(\mathbb{A})$ models are calculated based on the above metrics, using $w_i = 0.20$ as weights for dimension 1 and a linear decay with $R = 4$ (i.e., $W_i \in [0.33, 0.27, 0.20, 0.13, 0.07]$) for dimension 2. A more judicious choice of $w_i$ and/or $W_i$ could improve the results in RQ2 and RQ3, however our empirical study aims to provide a lower bound on the expected performance.

| activity | source | definition |
|---|---|---|
| `legacy` | git blame | #lines of code contributed by C that still survive in $R_i$ |
| `authored` | git log | #commits authored by C since $R_{i-1}$ |
| `committed` | git log | #commits committed by C since $R_{i-1}$ |
| `attributed` | git log | #commits since $R_{i-1}$ for which C is credited in the commit message under " |
| `reviewed` | mailing list | #commits since $R_{i-1}$ for which C has written at least one code review email |

Table 4.2 Concrete activity measures used for our empirical study on the Linux kernel. Each activity is measured for a particular contributor C of a specific subsystem S in a given release $R_i$.

### 4.4.4  Git-related Activity Measures

To calculate the Git-related activity measures of Table 4.2, i.e., all measures excluding `reviewed`, we cloned the official Linux kernel git repository, then checked out the Git tag corresponding to each analyzed release.

Bread-and-butter analysis of the Git log commits in the time span since the previous official release yields `authored` and `committed`, while simple regular expressions of the commit messages in the same logs obtains `attributed`. Finally, a standard git blame command yields, for each code line in the release under analysis, the last person touching it. This information allows to calculate `legacy`.

In kernel development, tags like "Signed-off-by", "Reviewed-by" and "Acked-by" are used as "a loose indication of review, so [...] need to be regarded as approximations only" **?**. Despite the warning of Corbet et al., `attributed` information is straightforward to obtain from commit messages, which is why we included this measure to complement the more strictly defined `reviewed` measure (calculated from reviewing data).

The next step is to lift up each contributor's Git-related activity measures to the subsystem-level, leveraging the file path information for each subsystem in the `MAINTAINERS` file. To do this, we identify for each commit the changed files, then map the commit to the subsystem(s) to which these files belong and aggregate the file-level measures to the subsystem level, for each contributor.

### 4.4.5  Linking Commits to Review Emails

In contrast to the developer `attributed` data obtained from the Git repository, the `reviewed` metric considers a second repository, i.e., the review environment. For the Linux kernel, code reviews are performed through mailing list discussions **???**. Patches are sent to one or more

of the various linux mailing lists (typically one per kernel subsystem), where anyone can step up and provide review comments simply by replying to the patch email. As such, the review comments of a specific patch are spread across one or more email threads.

Jiang et al. **??** have introduced a number of heuristics to link an accepted patch stored as a commit in the official Git repository to the (different versions of the) reviewed patch in the mailing list. We adopted the best performing heuristic of Jiang et al., which uses simple set intersection of the changed code lines between each email patch and Git commit. The heuristic matches a given Git commit C to the email patch P with which the change code line intersection is the largest and exceeds a given threshold of 4%. All emails in P's email thread are said to correspond to the review comments on P (and hence C).

To improve the line-based heuristic of Jiang et al., we have combined it with other heuristics. First of all, we observed that more and more kernel developers are using the commit message summary as the subject of their email threads. This summary is recommended to be between 50 and 72 characters [2] and appears before the body of a commit message. Hence, before applying the line-based matching of Jiang et al., we first check if there is a unique email patch P with subject identical to a commit C's commit summary. If so, we consider P and C to be a match, and do not need to run the more complex line-based matching algorithm.

If there is no such P, or multiple patches P have been found, we extract for each remaining commit the *author* and the *changed files*. We do the same for each remaining email patch. This information is then used to narrow down the search space of the line-based matching, by trying to match a commit only to email patches authored by the same developer and/or touching the same files. This substantially speeds up the matching process.

The remaining commits, i.e., the commits still not matched to a review, introduce noise to our measures. The reason for not finding a code review could be due to the reviews being sent to a mailing list that we did not analyze, or not being reviewed at all. We were granted read access to the database behind the Patchwork mailing list archive hosted by linux.org [3]. This Patchwork instance has been tracking 69 different Linux mailing lists since 2009, providing us with about 1.4 million patches. However, patches that were submitted through untracked mailing lists are not in our dataset, which explains the variability of matched commits across subsystems. Alternatively, the code change in the accepted commit could also have undergone substantial changes compared to the reviewed commit, for example due to rebasing, cherry-picking or squashing **?**.

Figure 4.1 shows the total percentage of matched commits from 2009 to the time of writing

---

2. `https://medium.com/@preslavrachev/what-s-with-the-50-72-rule-8a906f61f09c`
3. `https://patchwork.kernel.org/`

this paper, across the largest subdirectories of the kernel. It shows that this percentage varies greatly among the different subdirectories, with a minimum of 25% for the "net" subdirectory. This is why, in subsection 4.4.2, we filter out those subsystems whose average percentage of matched commits across the studied releases is lower than 60%. Note that we do not show the percentage of unmatched *patches* (only unmatched *commits*), since the unmatched email patches include those patches that were rejected during code review, and hence never showed up in the Git repository.

## 4.5   Case Study Results

This section discusses for each research question its motivation, specific approach and results.

### RQ1. How does expertise evolve in an open source project?

**Motivation :**

Open source software maintainers are responsible for the health of their subsystem. For example, each Linux kernel maintainer manages the changes proposed by developers to the subsystem they are responsible for, and shepherd those changes upstream towards the Git repository of Linus Torvalds (i.e., the official Linux kernel repository). Hence, their presence is vital to the kernel community.

Unfortunately, due to the unpredictable nature of life in general and open source software development in particular **??**, maintainers, for various reasons, one day will be forced to give up their responsibilities. In most cases, this means that another developer will have to take over the responsibility of maintainer.

Hence, this research question aims to analyze how often maintainership changes in kernel development. Furthermore, we are interested in understanding how much of the code base of official releases is "owned" by the subsystem maintainers, i.e., was originally developed by a maintainer. Since "git blame" is a popular means for finding experts **?**, our results will help us understand to what extent such a measure is reliable to measure expertise.

**Approach :**

To confirm the presence of changes in maintainership during the evolution of the Linux kernel, we analyzed the maintainers recorded in the `MAINTAINERS` file of releases *v3.5* to *v4.11* to identify how often maintainers (dis)appeared. Furthermore, for each studied release, we measure and plot these maintainers' `legacy`, which corresponds to the number of surviving code lines of a maintainer, as given by "git blame". We then validated the statistical signifi-

Figure 4.1 Percentage of matched commits in Linux subdirectories, from 2009 to the time of writing this paper.

cance of the change in `legacy` distribution between the first and last analyzed release using a Wilcoxon paired test. In case of a significant test result, we also provide the Cliff Delta effect size . An effect size smaller than 0.147 is deemed a "negligible" difference, smaller than 0.33 a "small" difference, smaller than 0.474 a "medium" difference and otherwise a "large" difference.

**Results :**

**23% of the studied subsystems saw changes in maintainership over the last 5 years.** Out of the 734 subsystems studied for RQ1, we counted 168 subsystems that experienced some sort of maintainership change. We counted 100 maintainer arrivals, 63 departures, and 88 replacements. These numbers confirm that maintainership change is common, even in mature open source systems like the Linux kernel. Furthermore, the median percentage of developers who are maintainers in the analyzed subsystems is 0.50% (mean of 0.90%), indicating that it is not straightforward to guess the next maintainer. These observations strengthen our case for more advanced expertise measures.

**The median maintainer `legacy` significantly decreases over time.** Figure 4.2 shows the evolution of the median percentage of maintainer `legacy` across all subsystems in each studied release. The plot shows a clear, steady decrease of this measure across releases in terms of median and variance. We confirm the significance of this decrease with a Wilcoxon paired test ($\alpha = 0.01$) between the first and last studied version, which yields a p-value of 2.2e-16.

Although the Cliff Delta value of 0.07 indicates only a negligible difference, this decreasing trend suggests that, if one limits the measure of expertise to the amount of surviving code originally authored by a maintainer, as was done by earlier work **?**, the expertise of maintainers globally seems to be decreasing over time. The next two research questions evaluate the use of a wider range of activity measures, across a range of releases, to obtain a more accurate measure of expertise.

**RQ2. How well does dimension 1 explain expertise?**

**Motivation :**

Prior work on expertise measures **?????** primarily are based on code activity, which can be

Figure 4.2 Median maintainer `legacy` across releases.

defined in terms of `legacy` and `committed`. As motivated in subsection 4.3.3, we believe that these two metrics do not capture the full breadth of contributor expertise activities. Indeed, the results in RQ1 indicate that the `legacy` of long-standing maintainers crumbles over time. Unless one assumes that this reflects a real drop in expertise over time, the only explanation is that existing experts reorient their focus to other activities, such as code review and email communication. Hence, this research question evaluates whether considering such additional activities is able to improve the identification of experts.

**Approach :**

To validate the ability of the measures in Table 4.2 to explain expertise, we evaluate how well the $footprint_j(\mathbb{A})$ measure involving those measures is able to identify the maintainers of Linux kernel subsystems. Those maintainers are the experts listed in the `MAINTAINERS` file of a kernel release.

For a given release and subsystem, we should find the maintainers in the top positions when ranked based on footprint values. The combination of activity measures $\mathbb{A}$ that is able to systematically yield the correct maintainers across subsystems and releases could be assumed to be a better indication of expertise.

In particular, we calculate two performance metrics :

$POS_N$  Percentage Of Subsystems for which at least one maintainer was ranked in the top N expert candidates

$POM_N$  Percentage Of Maintainers in the *whole* project who were ranked in the top N expert candidates for their subsystem

For these performance metrics, N is a threshold that can be varied. Our case study uses thresholds ranging from 1 to 5. It is important to note that, if the number of maintainers of a subsystem is larger than N, $POM_N$ could be penalized. To avoid this, we slightly changed the definition of $POM_N$ to be calculated only for the maintainers of all subsystems with at most N maintainers, instead of for all maintainers of the whole project. For example, $POM_1$ measures the percentage of top-recommended maintainers of subsystems with at most 1 maintainer.

To structure our analysis, we analyzed the performance of expertise models involving only one metric of Table 4.2, then analyzed models involving all combinations of `legacy` with one of the other 4 measures, and finally one model with all measures combined. We focus

explicitly on models involving `legacy` because it is a commonly used measure **?**, and hence we use its performance as baseline.

Since, for a given release, there is one $POS_N$ value and one $POM_N$ value, we calculate these metrics for each release, then study their distribution across the analyzed releases using boxplots. Figure 4.3 and Figure 4.4 show for each analyzed 𝔸 the distributions of $POS_N$ for N=1 and N=3, respectively, across the 22 studied Linux releases, while Figure 4.5 and Figure 4.6 the distributions of $POM_N$ for N=1 and N=3, respectively. We only show the plots for N=1 and N=3, as for higher values of N the plots remain more or less stable.

**Results :**

**`attributed` is the only single-measure model able to keep up with the multi-measure models.**

The results in Figure 4.3 and Figure 4.5 indicate that the first two individual measures, i.e., `legacy` and `authored`, are bad indicators of expertise compared to the other studied metrics. For example, in Figure 4.3, `legacy` only reaches a median $POS_N$ value of 47.22%, while `attributed` reaches a median $POS_N$ percentage of 58.4%.

**The models combining `legacy` with `committed`, `attributed` and/or `reviewed` perform the best.** Figure 4.3 shows indeed how only these four models read median percentages of 69.9%, while the other multi-measure models, especially the one involving only `legacy` and `authored`, are not able to outperform the best individual measure models.

These findings confirm the intuition that maintainers shifted focus from doing development (`authored`) themselves to mentoring others by controlling access to their subsystem's Git repository through committing and/or reviewing. As such, an expertise model only involving their own development (i.e., `legacy` and `authored`) is unable to explain the current kernel maintainers' expertise. In other words, modern expertise models should take into account the time spent reviewing code and pushing changes upstream.

$POS_N$ **increases to a median of 87.5% for larger N, with multi-measure models outperforming single-measure models by at least 17%.** Comparing Figure 4.4 to Figure 4.3 shows how the top multi-measure models for N=1 are able to increase their distance compared to even the best single-measure models (`attributed`). This, compared with a change in best performing single-measure models, indicates that a larger diversity in activity measures enables better identification of the two additional candidate maintainers. Indeed, by considering top performing contributors across a wider range of activities, there is a larger chance at least one real maintainer is found. Although the percentages of Figure 4.5 and Figure 4.6 cannot be compared directly to each other (cf. modified definition of $POM_N$), Fi-

Figure 4.3 Distribution of $POS_N$ for each combination of activity measures, for ranking threshold N = 1.

Figure 4.4 Distribution of $POS_N$ for each combination of activity measures, for ranking threshold N = 3.

Figure 4.5 Distribution of $POM_N$ for each combination of activity measures, for ranking threshold N = 1. These boxplots only consider subsystems with at most 1 maintainer.

Figure 4.6 Distribution of $POM_N$ for each combination of activity measures, for ranking threshold N = 3. These boxplots only consider subsystems with at most 3 maintainers.

gure 4.6 (for $POM_N$) shows a similar ranking of models as Figure 4.3 (for $POS_N$), confirming the findings for $POS_N$.

Table 4.3 shows the p-value and effect size of the Wilcoxon test between `attributed` and `legacy + committed` for figures 4.3, 4.4, 4.5, and 4.6. Each effect size being close to 1, we notice a **large** performance increase between `attributed` and `legacy + committed`.

Interesting to note is that, across all analyzed releases, the boxplots show a remarkable small variance, especially for N=3. Although this is partially due to the fact that less than 25% of the subsystems saw at least one maintainer change, it also indicates that our measures are stable across changes in the 5 activity measures used.

**RQ3. How well does dimension 2 explain expertise?**

**Motivation :**

The metrics analyzed in RQ2 reveal that traditional expertise metrics based solely on a contributor's own development productivity are not well suited to identify maintainers. Expertise models exploiting only the information available for the release under study, are able to obtain median $POS_N$ performance of up to 75% (N=1) and 90% (N=3).

However, we believe that adding a historical dimension considering also the activity in the last R releases would assist the model in two ways. On the one hand, long standing kernel developers' contributions should carry more weight than newcomers' contributions. On the other hand, analyzing data on multiple releases would control for cases where contributors' productivity was lower due to a variety of reasons, such as illness, vacation or work on other

| Type | Measure | N = 1 | N = 3 |
|------|---------|-------|-------|
| $POS_N$ | P-value | 4.27e-05 | 4.28e-05 |
| $POS_N$ | Cliff's delta | 0.99 | 1.0 |
| $POM_N$ | P-value | 4.27e-05 | 4.77e-07 |
| $POM_N$ | Cliff's delta | 0.84 | 1.0 |

Table 4.3 P-values and Cliff's delta values for the Wilcoxon paired tests ($\alpha = 0.01$) between `attributed` and `legacy + committed` for ranking thresholds N=1 and N=3 and for $POS_N$ and $POM_N$.

projects.

**Approach :**

For each studied kernel release, we calculated $footprint_j^R(\mathbb{A})$ for R=4, since this covers a time span of 60 to 70 days. For example, when looking at experts in release *v4.11*, we need to take into account data found for releases *v4.7, v4.8, v4.9, v4.10, and v4.11*. We repeat such analysis for each of the 22 releases. In this paper, we use linearly decaying weights $W_i$ to combine the individual $footprint_j(\mathbb{A})$ values across the five considered releases, since this scheme is less extreme than the exponential and logarithmic ones.

Similar to RQ2, we then use the footprint values to create, for each subsystem and release, a ranking of all contributors active in the five considered releases. We also use the same performance metrics as for RQ2, which allows us to compare the results of RQ3 to those of RQ2 to validate whether the historical dimension improves the model.

To save space, and since we found that, similar to RQ2, the combination of `legacy` and `committed` performs the best, we only show the results for this model (the rest of the data will be made available after the double-blind review). In particular, Figure 4.7 and Figure 4.8 show the $POS_N$ performance of the combined `legacy+committed` model without and with the history dimension, for ranking thresholds N ranging from 1 to 5. Figure 4.9 and Figure 4.10 show the corresponding $POM_N$ results.

Table 4.4 contains the results of Wilcoxon paired tests between the $POS_N$ values without and with history, for each N, and (similarly) between the $POM_N$ values without and with history, for each N. For each test, we also provide the Cliff Delta effect size .

**Results :**

**The history-aware `legacy+committed` footprint models perform significantly better than the history-unaware models.** Figure 4.7 and Figure 4.8 show how, except for N=3, the median performance of the history-aware expertise measure improves upon the history-unaware measure. If one considers only the first recommendation of the measure,

Figure 4.7 Distribution of $POS_N$ for the combined `legacy+committed` model (**without** history dimension), for different N.

Figure 4.8 Distribution of $POS_N$ for the combined `legacy+committed` model (**with** history dimension), for different N.

Figure 4.9 Distribution of $POM_N$ for the combined `legacy+committed` model (**without** history dimension), for different N. For each N, the boxplot only considers subsystems with at most N maintainers.

Figure 4.10 Distribution of $POM_N$ for the combined `legacy+committed` model (**with** history dimension), for different N. For each N, the boxplot only considers subsystems with at most N maintainers.

there is a median 73.6% chance that at least one maintainer is identified for a history-aware expertise model compared to 69.9% with the single-release model. This difference progressively decreases for higher N, which means that, for higher N, an expertise model considering `legacy+committed` on one release only is robust enough to assess expertise.

We find similar improvements for Figure 4.9 and Figure 4.10, except that the improvements due to history increase for larger N (and for N=1 there is no significant improvement). This is clearly shown by the p-values and effect sizes of the Wilcoxon paired tests in Table 4.4 ($alpha = 0.01$). As an effect size greater than 0.474 indicates a **large** increase in performance, we notice that 7 of the 8 significant differences have an effect size of at least 0.50.

## 4.6  Discussion

**Threats to validity :**

Threats to external validity prevent generalization of empirical results to other contexts. In particular, due to the abundant volume of data and presence of an oracle for expertise, our empirical evaluation only focused on 22 releases of the Linux kernel project. Hence, the study should be expanded to cover not only more kernel releases, but also other open (and closed) source projects. Furthermore, we considered only 5 expertise measures for our footprint models. Other measures, such as those mentioned in Table 4.1, should be studied to understand their impact on expertise.

Threats to construct validity involve risks regarding the measures used in the empirical

| Fig. | Measure | 1 | 2 | 3 | 4 | 5 |
|------|---------|---|---|---|---|---|
| 4.7/4.8 | P-value | 1.12e-4 | 8.76e-3 | 1.15e-2 | 1.57e-3 | 7.70e-4 |
| 4.7/4.8 | Cliff's delta | 0.62 | 0.50 | - | 0.55 | 0.51 |
| 4.9/4.10 | P-value | 1.27e-2 | 4.63e-3 | 5.13e-4 | 1.57e-5 | 1.88e-4 |
| 4.9/4.10 | Cliff's delta | - | 0.46 | 0.56 | 0.66 | 0.69 |

Table 4.4 P-values and Cliff's delta values for the Wilcoxon paired tests ($\alpha = 0.01$) between the $POS_N$ results of Figure 4.7 and Figure 4.8, and of Figure 4.9 and Figure 4.10, for ranking thresholds N=1 to N=5. A Cliff delta "-" indicates a non-significant test result, with p-value> $\alpha$

study. Of the five considered expertise measures, `reviewed` was the only one requiring noisy approximations. Except for cases where the email patch subject was identical to the Git commit message summary, there is a definite risk of false positive and false negative matches, as identified earlier by Jiang et al. **??**. This might explain the relatively weak performance of expertise models involving `reviewed`. However, no better alternatives exist for projects that use mailing lists for code review. Projects using web-based review environments like Gerrit do not have this issue, and will have perfect matching between commits and their reviews.

Finally, regarding threats to internal validity (i.e., confounding factors potentially explaining our findings), we mention the limited number of subsystems considered for RQ2 and RQ3. This number was the result of the data filtering in subsection 4.4.2 used to eliminate temporarily inactive subsystems. Furthermore, we used the `MAINTAINERS` file as oracle for expertise. Although this is the known reference in the Linux kernel community for finding the right maintainer to contact, this is a manually maintained text file that hence could contain inconsistencies (even though changes to it are peer-reviewed).

Finally, although maintainership is a form of expertise, there are other forms of expertise that our footprint models be indicators of that were not considered in our empirical study. As such, some of the false positive recommendations of our footprint rankings might actually be correct suggestions based on a different interpretation of expert, in which case our $POS_N$ and $POM_N$ results are lower bounds for the actual performance.

**Future work :**

Apart from addressing the threats to validity, other future work should consider different weights $w_i$ and $W_i$. The former weights consider different activities to be more relevant than others, while the latter weights would give more or less weights to older vs. newer releases. For example, comparison of exponential and logarithmic decaying weights to the linear decay used in our study could be interesting. Similarly, different values of $R$ for $footprint_j^R(\mathbb{A})$ should be evaluated.

Finally, whereas we used a top-down approach from expertise model to evaluation on an actual open source project, a bottom-up approach starting from the analysis of a project's or subsystem's maintainers before formulating expertise measures and models could provide complementary insights into different kinds expertise.

## 4.7 Conclusion

This paper argued about the need for expertise models considering a wide range of developer and other activities, and doing so across different snapshots of a project instead of just for one snapshot. Through an empirical study on 22 releases of the Linux kernel, we empirically showed how measures about an expert's own coding footprint (`legacy`) and her involvement in coordinating other project members (e.g., committing their commits and/or reviewing their code changes) significantly improves on coding-only expertise models. Furthermore, considering those measures across different releases significantly improved performance, with large effect size.

The simplest incarnation of our expertise model that software organizations should consider adopting involves (1) a developer's code `legacy` and number of changes `committed`, which are both readily obtainable from a Git log, calculated across (2) the last 5 releases. In future work, we will consider additional activity measures and empirically analyze other open source projects.

# CHAPTER 5 THIRD TOPIC WITH A VERY VERY LONG TITLE THAT TAKES MORE THAN ONE LINE

Lorem ipsum dolor sit amet, non faucibus ut, ante integer tristique odio vitae turpis in. Euismod ullamcorper urna eget sollicitudin consectetuer, dolor a. Ridiculus volutpat fusce, montes ipsum placerat, eu malesuada maecenas a odio per, est pellentesque integer auctor sed ut sed, lectus sodales orci ornare. Donec neque turpis vehicula. Duis vel sapien nec massa lobortis nonummy. Feugiat ultrices urna mauris.

Potenti erat molestie ridiculus placerat, viverra ut felis porttitor, rhoncus accumsan non, dui magna quam justo, ultrices massa ut phasellus donec viverra mauris. Mauris a, dictumst risus a ornare velit nulla ultricies, neque leo pellentesque, sit sed et suscipit excepteur aenean. Venenatis sodales, odio nostra in id nobis scelerisque, venenatis sociosqu gravida blandit orci pellentesque, tincidunt velit sed elementum lacus pretium nunc, aenean vel dui id. Elit placerat id dui nunc mollis, diam sapien porta, ipsam elit magna imperdiet amet, erat feugiat, et eros morbi feugiat velit fringilla. Lacinia phasellus lacinia magna nunc sed, a rhoncus, sem eget, dui aliquam sit sed leo beatae non, quisque justo dignissim.

Torquent curabitur magnis nullam viverra scelerisque, per lacus pellentesque vivamus, mauris aliquam sem lacus vivamus nullam porta. Vivamus donec maecenas nunc orci massa, orci neque luctus leo non, mauris quis metus sagittis. Voluptatibus gravida interdum. Magna duis nulla odio lacus fugiat non. Magna fusce nunc, eget pellentesque nec. Imperdiet non magna sollicitudin pellentesque, fusce erat interdum diam tellus vel, vitae iaculis lectus varius suspendisse. Ac vel a in semper tellus, lobortis sed, ipsum volutpat. Mauris a nunc aliquam metus nec, eu et id risus, diam integer molestie suspendisse, sed wisi. Metus sed justo sodales sapien molestie, suspendisse sem viverra ac proin, lorem luctus at tellus, velit mi morbi orci in vestibulum, dignissim urna ornare id donec. Suspendisse non enim euismod odio elit mauris, consectetuer pellentesque faucibus velit ante lacinia sed.

Et dui erat. Wisi lorem eleifend cursus do donec, sed vel fermentum nec, a a in pharetra. Ultricies risus, eget habitasse in, consectetuer metus in auctor ac pellentesque curabitur, pulvinar aliquet eget. Mattis eget venenatis dolor, nunc sem sed massa, urna scelerisque a magnis, neque elit nec aliquam nonummy ac accusantium. Id vivamus nunc, erat justo tellus, scelerisque habitasse accumsan tellus, pede sem vestibulum velit in et eleifend. Nulla massa aenean integer dui. Suscipit nunc purus, rutrum velit, mi torquent elementum in tincidunt. Maecenas nulla integer fringilla dapibus tellus sit, enim amet magna eu erat, libero consectetuer nisl sapien, in ultricies neque arcu sodales sagittis.

Lorem ipsum dolor sit amet, non faucibus ut, ante integer tristique odio vitae turpis in. Euismod ullamcorper urna eget sollicitudin consectetuer, dolor a. Ridiculus volutpat fusce, montes ipsum placerat, eu malesuada maecenas a odio per, est pellentesque integer auctor sed ut sed, lectus sodales orci ornare. Donec neque turpis vehicula. Duis vel sapien nec massa lobortis nonummy. Feugiat ultrices urna mauris.

Potenti erat molestie ridiculus placerat, viverra ut felis porttitor, rhoncus accumsan non, dui magna quam justo, ultrices massa ut phasellus donec viverra mauris. Mauris a, dictumst risus a ornare velit nulla ultricies, neque leo pellentesque, sit sed et suscipit excepteur aenean. Venenatis sodales, odio nostra in id nobis scelerisque, venenatis sociosqu gravida blandit orci pellentesque, tincidunt velit sed elementum lacus pretium nunc, aenean vel dui id. Elit placerat id dui nunc mollis, diam sapien porta, ipsam elit magna imperdiet amet, erat feugiat, et eros morbi feugiat velit fringilla. Lacinia phasellus lacinia magna nunc sed, a rhoncus, sem eget, dui aliquam sit sed leo beatae non, quisque justo dignissim.

Torquent curabitur magnis nullam viverra scelerisque, per lacus pellentesque vivamus, mauris aliquam sem lacus vivamus nullam porta. Vivamus donec maecenas nunc orci massa, orci neque luctus leo non, mauris quis metus sagittis. Voluptatibus gravida interdum. Magna duis nulla odio lacus fugiat non. Magna fusce nunc, eget pellentesque nec. Imperdiet non magna sollicitudin pellentesque, fusce erat interdum diam tellus vel, vitae iaculis lectus varius suspendisse. Ac vel a in semper tellus, lobortis sed, ipsum volutpat. Mauris a nunc aliquam metus nec, eu et id risus, diam integer molestie suspendisse, sed wisi. Metus sed justo sodales sapien molestie, suspendisse sem viverra ac proin, lorem luctus at tellus, velit mi morbi orci in vestibulum, dignissim urna ornare id donec. Suspendisse non enim euismod odio elit mauris, consectetuer pellentesque faucibus velit ante lacinia sed.

Et dui erat. Wisi lorem eleifend cursus do donec, sed vel fermentum nec, a a in pharetra. Ultricies risus, eget habitasse in, consectetuer metus in auctor ac pellentesque curabitur, pulvinar aliquet eget. Mattis eget venenatis dolor, nunc sem sed massa, urna scelerisque a magnis, neque elit nec aliquam nonummy ac accusantium. Id vivamus nunc, erat justo tellus, scelerisque habitasse accumsan tellus, pede sem vestibulum velit in et eleifend. Nulla massa aenean integer dui. Suscipit nunc purus, rutrum velit, mi torquent elementum in tincidunt. Maecenas nulla integer fringilla dapibus tellus sit, enim amet magna eu erat, libero consectetuer nisl sapien, in ultricies neque arcu sodales sagittis.

# CHAPTER 6    CONCLUSION

Text.

## 6.1    Advancement of knowledge

Text.

## 6.2    Limits and constraints

Text.

## 6.3    Recommendations

Text.

# REFERENCES

C. Abondo, "Gestion de la Qualité de Service dans les systèmes mobiles de prochaine génération", Thèse de doctorat, École Polytechnique de Montréal, Juin 2005.

T. Bao et X. Zhang, "On-the-fly detection of instability problems in floating-point program execution", *SIGPLAN Not.*, vol. 48, no. 10, pp. 817–832, Oct. 2013. DOI : `10.1145/2544173.2509526`. En ligne : `http://doi.acm.org/10.1145/2544173.2509526`

F. Benz, A. Hildebrandt, et S. Hack, "A dynamic program analysis to find floating-point accuracy problems", dans *ACM SIGPLAN Notices*, vol. 47, Juin 2012, pp. 453 – 462.

A. Dabrowski, P. Pawowski, M. Stankiewicz, et F. Misiorek, "Quasi-maximum accuracy floating-point computations with GPGPU for applications in digital signal processing", dans *SPA 2011 - Signal Processing : Algorithms, Architectures, Arrangements, and Applications - Conference Proceedings*, Poznan, Poland, 2011, pp. 144 – 148.

F. De Dinechin, C. Lauter, et G. Melquiond, "Certifying the floating-point implementation of an elementary function using Gappa", *IEEE Transactions on Computers*, vol. 60, no. 2, pp. 242 – 253, 2011. DOI : `10.1109/TC.2010.128`

L. de Figueiredo et J. Stolfi, "Affine arithmetic : Concepts and applications", *Numerical Algorithms*, vol. 37, no. 1-4, pp. 147–158, 2004. DOI : `10.1023/B:NUMA.0000049462.70970.b6`

V. D'Silva, L. Haller, D. Kroening, et M. Tautschnig, "Numeric bounds analysis with conflict-driven learning", dans *Tools and Algorithms for the Construction and Analysis of Systems*, série Lecture Notes in Computer Science, C. Flanagan et B. König, éds. Springer Berlin Heidelberg, 2012, vol. 7214, pp. 48–63. DOI : `10.1007/978-3-642-28756-5_5`

J.-C. Filliâtre et C. Marché, "The Why/Krakatoa/Caduceus platform for deductive program verification", dans *Computer Aided Verification*, série Lecture Notes in Computer Science, W. Damm et H. Hermanns, éds. Springer Berlin Heidelberg, 2007, vol. 4590, pp. 173–177. DOI : `10.1007/978-3-540-73368-3_21`

L. Fousse, G. Hanrot, V. Lefèvre, P. Pélissier, et P. Zimmermann, "MPFR : A multiple-precision binary floating-point library with correct rounding", *ACM Transactions on Mathematical Software*, vol. 33, no. 2, pp. 1–14, Juin 2007. DOI : `10.1145/1236463.1236468`

D. Goldberg, "What every computer scientist should know about floating-point arithmetic", *Computing Surveys*, vol. 23, no. 1, pp. 5–48, Mars 1991. En ligne : `http://docs.oracle.com/cd/E19957-01/806-3568/ncg_goldberg.html`

E. Goubault, "Static analysis of the precision of floating-point operations", dans *Static Analysis. 8th International Symposium, SAS 2001. Proceedings*, Berlin, Germany, 2001, pp. 234 – 59.

E. Goubault, S. Putot, P. Baufreton, et J. Gassino, "Static analysis of the accuracy in control systems : Principles and experiments", dans *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 4916 LNCS, Berlin, Germany, 2008, pp. 3 – 20. DOI : `10.1007/978-3-540-79707-4_3`

D. W. Harder et R. Khoury, *Numerical Analysis for Engineering*, University of Waterloo, 200 University Avenue West, Waterloo, Ontario, Canada N2L 3G1, 2005. En ligne : `https://ece.uwaterloo.ca/~dwharder/NumericalAnalysis/`

N. J. Higham, *Accuracy and Stability of Numerical Algorithms.* SIAM, 2002.

F. Johansson *et al.*, *mpmath : a Python library for arbitrary-precision floating-point arithmetic (version 0.18)*, Déc. 2013. En ligne : `http://mpmath.org/`

E. Nichols, L. J. McDaid, et N. H. Siddique, "Case study on a self-organizing spiking neural network for robot navigation", *International Journal of Neural Systems*, vol. 20, no. 6, pp. 501–508, 2010.

NVIDIA, *CUDA C Programming Guide*, 5e éd., NVIDIA, 2012. En ligne : `http://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html`

J. Postel, *RFC791 : Internet Protocol*, IETF, Sep. 1981.

A. Tanenbaum, *Computer Networks (fourth edition).* Prentice-Hall International, Inc., 2002.

C. Tseng, G. Lee, R. Liu, et T. Wang, "HMRSVP : a Hierarchical Mobile RSVP Protocol", *Wireless Networks*, vol. 9, no. 2, pp. 95–102, Mars 2003.

N. Whitehead et A. Fit-Florea, "Precision & performance : Floating point and IEEE 754 compliance for NVIDIA GPUs", 2011. En ligne : `https://developer.nvidia.com/sites/default/files/akamai/cuda/files/NVIDIA-CUDA-Floating-Point.pdf`

L. Zhang, S. Pierre, et L. Marchand, "Optimization of Handover Performance for FMIPv6", dans *Intelligence in Communication Systems*, 2005, pp. 169–178.

# ANNEXE A    DEMO

Texte de l'annexe A. Remarquez que la phrase précédente se termine par une lettre majuscule suivie d'un point. On indique explicitement cette situation à LaTeX afin que ce dernier ajuste correctement l'espacement entre le point final de la phrase et le début de la phrase suivante.

# ANNEXE B    ANOTHER APPENDIX

Text in «landscape» mode.

## ANNEXE C    ONE MORE APPENDIX

Text.