



**HAL**  
open science

## À la racine du parallélisme

Thomas Bonald, Céline Comte, Fabien Mathieu

► **To cite this version:**

Thomas Bonald, Céline Comte, Fabien Mathieu. À la racine du parallélisme. [Rapport de recherche] Telecom ParisTech. 2017. hal-01476889v1

**HAL Id: hal-01476889**

**<https://inria.hal.science/hal-01476889v1>**

Submitted on 26 Feb 2017 (v1), last revised 2 May 2017 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# À la racine du parallélisme

Thomas Bonald<sup>1</sup>, Céline Comte<sup>2</sup><sup>1</sup> et Fabien Mathieu<sup>2</sup> †

<sup>1</sup>Télécom ParisTech, Université Paris-Saclay, France

<sup>2</sup>Nokia Bell Labs, France

---

Nous considérons un clusters de serveurs traitant des requêtes en parallèle. Si les clients ont en général intérêt à ce que leurs requêtes soient traitées par le plus grand nombre de serveurs, l'impact du parallélisme sur les serveurs est moins clair : trop faible, il ne permet pas d'utiliser pleinement les ressources disponibles ; trop fort, il risque d'encombrer inutilement les serveurs de requêtes en attente. Nous étudions ce phénomène à l'aide d'un modèle de files d'attente où les requêtes arrivent selon un processus de Poisson et requièrent des traitements dont le volume suit une loi exponentielle. Chaque nouvelle requête est affectée à un certain nombre de serveurs, choisis de manière aléatoire, uniforme, et indépendante de l'état du système. Chaque serveur traite ses requêtes dans leur ordre d'arrivée. Nous montrons qu'il existe un degré de parallélisme qui minimise le nombre moyen de requêtes présentes dans chaque serveur. Ce degré optimal est de l'ordre de la racine carrée du nombre de serveurs pour une charge faible à modérée, et décroît jusqu'à deux à très forte charge.

**Mots-clefs :** Cluster de serveurs, répartition de charge, parallélisme.

---

## 1 Introduction

Paralléliser les tâches dans les centres de calcul permet de réduire leur temps de traitement et d'optimiser l'utilisation des ressources. Même en l'absence de coût de parallélisation, par exemple lié à la communication entre les serveurs, il doit exister un degré de parallélisme optimal dans un tel système : trop faible, il ne permet pas d'utiliser pleinement les ressources disponibles ; trop fort, il risque d'encombrer inutilement les serveurs de requêtes. Les modèles de files d'attente traditionnels ne permettent pas de capturer ce phénomène [Kle75]. Nous nous intéressons ici à un modèle plus récent introduit par Gardner et al. [GHBSW<sup>+</sup>17] dans lequel les requêtes peuvent être parallélisées sur un nombre fixe de serveurs choisis de façon aléatoire, uniforme et indépendante de l'état du système, chaque serveur traitant ses requêtes dans leur ordre d'arrivée. Nous montrons qu'il existe un degré de parallélisme pour lequel le nombre moyen de requêtes sur chaque serveur est minimal. De manière assez surprenante, ce degré optimal est de l'ordre de la racine du nombre de serveurs à charge faible à modérée et décroît jusqu'à 2 lorsque la charge est très forte.

## 2 Modèle

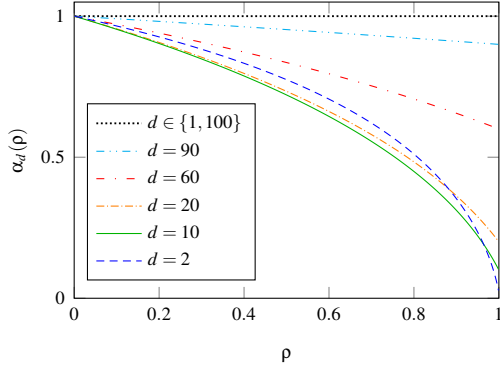
Soit un cluster de  $N$  serveurs, chacun ayant  $\mu > 0$  pour capacité de traitement. Des requêtes arrivent dans le système selon un processus de Poisson d'intensité  $N\lambda$  ; chaque requête est répartie sur  $d$  serveurs choisis au hasard de façon indépendante de l'état du système. La quantité de travail de chaque requête suit une loi exponentielle de moyenne unitaire. Les requêtes quittent le système à la fin de leur service.

Chaque serveur traite les requêtes qui lui ont été attribuées séquentiellement par ordre d'arrivée. Une requête peut être traitée efficacement en parallèle par plusieurs serveurs, de sorte que son taux de service est la somme des capacités des serveurs en train de le servir. Notons que ce taux peut augmenter suite au départ de requêtes précédentes. La charge totale de la file est notée  $\rho = \frac{N\lambda}{N\mu} = \frac{\lambda}{\mu}$ . Comme tous les serveurs sont interchangeables,  $\rho$  donne aussi la charge de chaque serveur.

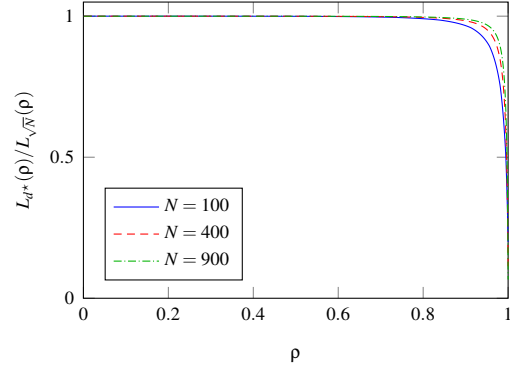
Ce système peut être vu comme une file d'attente multi-classe multi-serveur, chaque classe de requêtes correspondant à l'une des  $\binom{N}{d}$  attributions possibles. Sans grande surprise, la file est stable si et seulement

---

†Les auteurs sont membres du LINCS, voir <http://www.lincs.fr>



**FIGURE 1:** Nombre moyen normalisé de requêtes à traiter par chaque serveur en fonction de la charge  $\rho$  ( $N = 100$ ).



**FIGURE 2:** Optimalité du choix  $d = \sqrt{N}$  en fonction de la charge  $\rho$ .

si  $\rho < 1$  (cf [GHBSW<sup>+</sup>17]), ce que nous supposons vérifié dans la suite. La file est décrite par un état agrégé  $(x, n)$ , où  $x$  est le nombre de requêtes présentes et  $n$  le nombre de serveurs actifs (c'est-à-dire en train de traiter une requête). Cet état définit un processus stochastique, en général non markovien, sur

$$\mathcal{S} = \{(0, 0)\} \cup \{(x, n) \in \mathbb{N}^2 : x \geq 1 \text{ et } d \leq n \leq \min(N, dx)\}.$$

On note  $\pi$  sa distribution à l'état stationnaire. Gardner *et al.* ont prouvé dans [GHBSW<sup>+</sup>17, Théorème 4] que  $\pi$  satisfait la récurrence suivante : pour chaque  $(x, n) \in \mathcal{S} \setminus \{(0, 0)\}$ , on a

$$\pi(x, n) = \frac{N\rho}{n} \sum_{\ell=0}^d \frac{\binom{n-\ell}{d-\ell} \binom{N-n+\ell}{\ell}}{\binom{N}{d}} \pi(x-1, n-\ell), \quad \text{avec la convention } \pi(x, n) = 0 \text{ si } (x, n) \notin \mathcal{S}. \quad (1)$$

Par la suite, on note  $(\mathbf{X}, \mathbf{N})$  un couple de variables aléatoires suivant cette loi de probabilité.

### 3 Taille moyenne de la file d'attente par serveur

La formule suivante a été prouvée dans [GHBSW<sup>+</sup>17, Théorème 1] :

$$\mathbb{E}(\mathbf{X}) = \sum_{n=d}^N \frac{\rho}{\frac{\binom{N-1}{d-1}}{\binom{n-1}{d-1}} - \rho}. \quad (2)$$

On note  $L_d(\rho)$  le nombre moyen de requêtes sur chaque serveur. Tous les serveurs sont interchangeables et chaque nouvelle requête est répliquée sur  $d$  serveurs distincts, de sorte que

$$L_d(\rho) = \frac{d\mathbb{E}(\mathbf{X})}{N} = \frac{d}{N} \sum_{n=d}^N \frac{\rho}{\frac{\binom{N-1}{d-1}}{\binom{n-1}{d-1}} - \rho}. \quad (3)$$

Dans les cas particuliers où  $d = 1$  et  $d = N$ , on obtient le nombre moyen  $\rho/(1 - \rho)$  de requêtes dans une file d'attente  $M/M/1$  sous la charge  $\rho$ . Lorsque l'on regarde l'évolution de la file d'attente de chaque serveur, les deux situations sont effectivement équivalentes à des échelles de temps différentes. Dans le premier cas, on a  $N$  files  $M/M/1$  indépendantes : chaque serveur voit arriver les requêtes au taux  $\lambda$  et les traite à vitesse  $\mu$  selon la discipline FIFO. Dans le second cas, on a une unique file  $M/M/1$  : le taux d'arrivée des requêtes à chaque serveur est  $N\lambda$  et celles-ci sont servies par ordre d'arrivée à vitesse  $N\mu$ . Il est montré dans l'annexe A que ces deux configurations constituent un pire cas pour le nombre moyen de requêtes par serveur.

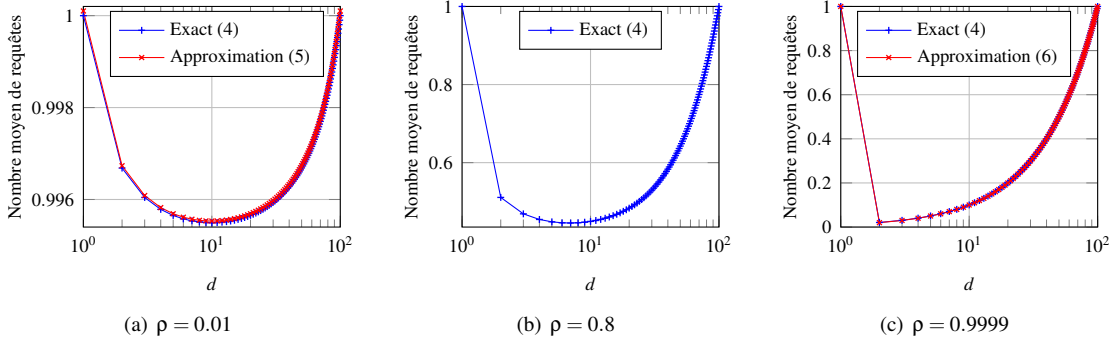


FIGURE 3:  $\alpha_d(\rho)$  en fonction de  $d$ , pour différentes valeurs de  $\rho$  ( $N = 100$ ).

Notre objectif est de minimiser l'encombrement moyen  $L_d(\rho)$  et de quantifier le gain par rapport au pire cas  $d = 1$ , sans parallélisme. On s'intéresse donc à la quantité suivante, comprise entre 0 et 1 :

$$\alpha_d(\rho) \equiv \frac{L_d(\rho)}{L_1(\rho)} = \frac{\frac{d\mathbb{E}(\mathbf{X})}{N}}{\frac{\rho}{1-\rho}} = \frac{d}{N} \sum_{n=d}^N \frac{1-\rho}{\binom{N-1}{d-1} - \rho}. \quad (4)$$

L'évolution de  $\alpha_d(\rho)$  en fonction de la charge  $\rho$  pour  $N = 100$  serveurs est représentée sur la Figure 1 pour différentes valeurs de  $d$ . On observe tout d'abord que, exceptées les valeurs extrêmes  $d = 1$  et  $d = 100$ , le parallélisme est d'autant plus efficace que la charge est forte. On voit également que le choix  $d = 10$  est optimal sauf pour une charge  $\rho$  proche de 1, où  $d = 2$  devient meilleur.

Afin de mieux comprendre le phénomène, la Figure 3 étudie trois « tranches » de la Figure 1, à savoir la valeur de  $\alpha_d(\rho)$  en fonction de  $d$  pour une charge  $\rho$  faible, forte et extrême. À charge faible (Figure 3(a)), on observe une certaine symétrie de la courbe, en échelle logarithmique pour  $d$ , qui amène à penser que le choix  $d = \sqrt{N}$  est optimal. Cette symétrie se brise lentement quand  $\rho$  augmente, et à forte charge (Figure 3(b)), le choix  $d = \sqrt{N}$  n'est plus optimal, même s'il reste assez efficace. À charge extrême (Figure 3(c)), la symétrie est complètement brisée : la courbe est croissante pour  $d \geq 2$ .

Pour montrer l'intérêt de choisir pour  $d$  la racine carrée du nombre de serveurs, nous avons calculé à  $N$  et  $\rho$  fixés la valeur  $d^*$  de  $d$  qui minimise  $L_d(\rho)$  et déduit l'optimalité relative  $L_{d^*}(\rho)/L_{\sqrt{N}}(\rho)$  du choix  $d = \sqrt{N}$ . Les résultats, reportés sur la Figure 2, confirment une quasi-optimalité de  $d = \sqrt{N}$  sauf à très forte charge. On peut également remarquer que la baisse d'optimalité semble retardée pour  $N$  grand.

Nous allons maintenant prouver l'optimalité de  $d = \sqrt{N}$  à faible charge et de  $d = 2$  à forte charge, la quasi-optimalité de  $d = \sqrt{N}$  à charge modérée restant pour l'instant une conjecture, même si elle est mise en évidence par les évaluations numériques.

## 4 Étude à faible charge

À faible charge, la performance du système est dictée par ce qu'il se passe en présence d'une ou deux requêtes. Celle qui est en tête de file est systématiquement en service sur  $d$  serveurs ; la seconde, si elle est présente, est en service sur  $d - k$  serveurs, où  $k$  est le nombre de *collisions*, c'est-à-dire de serveurs qui sont attribués aux deux requêtes. Il y a en moyenne  $d^2/N$  collisions, puisque chacun des  $d$  serveurs attribués à la seconde requête peut avoir déjà été attribué à la première avec probabilité  $d/N$ . Intuitivement, on comprend que la valeur  $d = \sqrt{N}$  marque une transition : si  $d \ll \sqrt{N}$ , il y a peu de collisions, mais un parallélisme faible, alors que si  $d \gg \sqrt{N}$ , l'effet du parallélisme est amoindri par un grand nombre de collisions. Le résultat suivant permet de formaliser cette situation.

**Proposition 1.** Pour  $1 \leq d \leq N$ , on a

$$\alpha_d(\rho) \simeq \frac{1-\rho}{1 - \frac{\rho}{2 - (\frac{1}{d} + \frac{1}{N})}} \text{ lorsque } \rho \rightarrow 0. \quad (5)$$

La résolution de (5) montre que la valeur de  $d$  qui minimise  $\alpha_d(\rho)$  est bien  $\sqrt{N}$ . La validité de l'approximation peut être observée sur la Figure 3(a) pour  $N = 100$  et  $\rho = 0.01$ . On donne ici une idée de la preuve de (5), la preuve complète étant disponible dans l'annexe B.

*Idée de la preuve.* On fait tendre  $N\rho/d$  vers 0. Quand  $N\rho/d$  est petit, on peut ignorer les états contenant plus de deux requêtes car leur impact dans  $\alpha_d(\rho)$ , en  $O((N\rho/d)^2)$ , est négligeable devant le reste : en développant puis en simplifiant l'expression (4) de  $\alpha_d(\rho)$  en fonction de  $\mathbb{E}(\mathbf{X})$ , on obtient

$$\alpha_d(\rho) = (1 - \rho) \frac{1 + \frac{N\rho}{d} \sum_{k=0}^d \frac{1}{1 - \frac{k}{2d}} p_k + O\left(\left(\frac{N\rho}{d}\right)^2\right)}{1 + \frac{N\rho}{d} + O\left(\left(\frac{N\rho}{d}\right)^2\right)},$$

où  $p_k = \binom{d}{k} \binom{N-d}{d-k} / \binom{N}{d}$ , pour  $k = 0, 1, \dots, d$ , est la loi hypergéométrique de paramètres  $N, d, d$  (probabilité de choisir  $k$  serveurs occupés parmi  $d$  serveurs choisis au hasard parmi  $N$  serveurs dont  $d$  sont occupés). Il reste à donner une approximation de la somme où intervient cette loi. On distingue trois régimes selon les valeurs de  $d$  :  $d$  petit devant  $\sqrt{N}$ ,  $d$  de l'ordre de  $\sqrt{N}$  et  $d$  proche de  $N$  (ou plus précisément,  $N - d$  petit devant  $\sqrt{N}$ ). Dans les trois cas, on montre que

$$\sum_{k=0}^d \frac{1}{1 - \frac{k}{2d}} p_k \simeq 1 + \frac{d}{N} \frac{1}{2 - \frac{1}{d} - \frac{d}{N}}, \text{ dont on déduit le résultat.}$$

□

## 5 Étude à forte charge

Le comportement à forte charge est décrit par la proposition suivante.

**Proposition 2.** *Pour  $2 \leq d \leq N$ , on a*

$$\alpha_d(\rho) \simeq \frac{d}{N} \text{ lorsque } \rho \rightarrow 1. \quad (6)$$

En se souvenant que  $\alpha_1(\rho) = 1$ , on en déduit que le degré de parallélisme qui minimise le nombre moyen de requêtes par serveur à forte charge est  $d = 2$ . La validité de (6) est illustrée sur la Figure 3(c) pour  $N = 100$  et  $\rho = 1 - 1/N^2 = 0.9999$ . L'idée de la preuve, donnée dans l'annexe C, repose sur un développement limité de  $\alpha_d(\rho)$  quand  $\rho$  tend vers 1. Dès que  $d \geq 2$ , le terme dominant dans (2) est  $\rho/(1 - \rho)$ , correspondant à  $n = N$ . Le nombre moyen de requêtes dans le cluster s'approche ainsi de  $\rho/(1 - \rho)$  à très forte charge. En revenant à  $\alpha_d(\rho)$ , on obtient le résultat.

## 6 Conclusion

Nous avons étudié un modèle de cluster de serveurs avec traitement des requêtes en parallèle et montré qu'il existe un degré de parallélisme qui minimise le nombre moyen de requêtes sur chaque serveur. Ce degré optimal est de l'ordre de la racine du nombre de serveurs à faible charge et décroît jusqu'à 2 à forte charge. Dans les travaux futurs, nous souhaitons prouver les intuitions données par les résultats numériques et étudier la sensibilité du degré optimal à la loi de la taille des requêtes.

## Références

- [GHBSW<sup>+</sup>17] K. Gardner, M. Harchol-Balter, A. Scheller-Wolf, M. Velednitsky, and S. Zbarsky. Redundancy-d : The power of d choices for redundancy. *Operations Research (to appear, available at <http://www.cs.cmu.edu/~harchol/Papers/OR16.pdf>)*, 2017.
- [Kle75] L. Kleinrock. *Queueing Systems, Volume I : Theory*. Wiley Interscience, New York, 1975.

À la racine du parallélisme

## Annexes

La notation  $f(N, d, \rho) = O(g(N, d, \rho))$  est utilisée dans la suite pour signifier qu'une fonction  $f$  de  $N$ ,  $d$  et  $\rho$  est dominée par une autre fonction  $g$  dans le sens où il existe une constante  $C > 0$  telle que

$$f(N, d, \rho) \leq Cg(N, d, \rho),$$

lorsque  $N$  et  $d$  tendent vers  $+\infty$  et  $\rho$  est proche de 0 ou de 1 pour les Annexes B et C respectivement. L'objectif est de présenter des résultats aussi généraux que possible, sans s'attacher à un régime particulier de passage à l'échelle. Dans la suite, on note aussi  $n \wedge m$  le minimum de deux entiers relatifs  $n$  et  $m$ , et  $n^+ = n \vee 0$  le maximum d'un entier relatif  $n$  et de 0.

### A Majoration de $L_d(\rho)$

Soit  $d \in \{1, \dots, N\}$ . On va montrer que le système se comporte au moins aussi bien qu'une file d'attente  $M/M/1$  de même charge, c'est-à-dire que

$$L_d(\rho) \leq \frac{\rho}{1-\rho}.$$

D'après (3), on a

$$L_d(\rho) = \frac{d}{N} \sum_{n=d}^N \frac{\rho}{\frac{\binom{n-1}{d-1}}{\binom{n-1}{d-1}} - \rho} = \frac{d}{N} \sum_{n=d}^N \frac{\binom{n-1}{d-1}}{\binom{n-1}{d-1}} \rho \frac{1}{1 - \frac{\binom{n-1}{d-1}}{\binom{n-1}{d-1}} \rho}.$$

Puisque  $\rho < 1$ , on peut utiliser le développement en série entière suivant :

$$L_d(\rho) = \frac{d}{N} \sum_{n=d}^N \frac{\binom{n-1}{d-1}}{\binom{n-1}{d-1}} \rho \sum_{m=0}^{+\infty} \left[ \frac{\binom{n-1}{d-1}}{\binom{n-1}{d-1}} \rho \right]^m = \frac{d}{N} \sum_{n=d}^N \sum_{m=1}^{+\infty} \left[ \frac{\binom{n-1}{d-1}}{\binom{n-1}{d-1}} \rho \right]^m = \sum_{m=1}^{+\infty} \left( \frac{d}{N} \sum_{n=d}^N \left[ \frac{\binom{n-1}{d-1}}{\binom{n-1}{d-1}} \right]^m \right) \rho^m.$$

Comme  $0 \leq \frac{\binom{n-1}{d-1}}{\binom{n-1}{d-1}} \leq 1$  et  $x^m \leq x$  si  $0 \leq x \leq 1$ , on obtient la majoration

$$L_d(\rho) \leq \sum_{m=1}^{+\infty} \left( \frac{d}{N} \sum_{n=d}^N \frac{\binom{n-1}{d-1}}{\binom{n-1}{d-1}} \right) \rho^m.$$

D'après l'identité de la crosse de hockey (*Hockey-stick identity*, dite aussi *Christmas stocking identity*), on a par ailleurs

$$\sum_{n=d}^N \binom{n-1}{d-1} = \sum_{n=d-1}^{N-1} \binom{n}{d-1} = \binom{(N-1)+1}{(d-1)+1} = \binom{N}{d},$$

d'où

$$\frac{d}{N} \sum_{n=d}^N \frac{\binom{n-1}{d-1}}{\binom{n-1}{d-1}} = \frac{d}{N} \frac{\sum_{n=d}^N \binom{n-1}{d-1}}{\binom{n-1}{d-1}} = \frac{d}{N} \frac{\binom{N}{d}}{\binom{N-1}{d-1}} = 1.$$

On conclut ainsi que

$$L_d(\rho) \leq \sum_{m=1}^{+\infty} \rho^m = \frac{\rho}{1-\rho}.$$

## B Preuve de la Proposition 1 de la Partie 4

On cherche à montrer que pour  $N\rho/d$  petit, on a

$$\alpha_d(\rho) \simeq \frac{1-\rho}{1-\frac{\rho}{2-(\frac{1}{d}+\frac{1}{N})}}. \quad (5)$$

Sans perte de généralité, on suppose pour toute la partie B que  $N\rho/d \leq 1/2$ . Comme annoncé dans l'esquisse de preuve de la Partie 4, on commence par montrer que les états contenant plus de deux requêtes peuvent être ignorés car leur impact dans  $\alpha_d(\rho)$ ,  $O((N\rho/d)^2)$ , est dominé par l'impact des états à deux requêtes ou moins. Ceci est formalisé dans le Lemme suivant et son Corollaire.

**Lemme 1.** *Pour tout  $x \geq 1$ ,*

$$\sum_{n=0}^N \pi(x, n) \leq \frac{N\rho}{d} \sum_{n=0}^N \pi(x-1, n).$$

*Démonstration.* Soit  $x \geq 1$ . (1) montre que  $\pi(x, n) = 0$  si  $n < d$ . On a donc, toujours d'après (1),

$$\sum_{n=0}^N \pi(x, n) \leq \frac{N\rho}{d} \sum_{n=d}^N \sum_{\ell=0}^d \frac{\binom{n-\ell}{d-\ell} \binom{N-n+\ell}{\ell}}{\binom{N}{d}} \pi(x-1, n-\ell).$$

Grâce à un changement de variable, on obtient

$$\begin{aligned} \sum_{n=d}^N \sum_{\ell=0}^d \frac{\binom{n-\ell}{d-\ell} \binom{N-n+\ell}{\ell}}{\binom{N}{d}} \pi(x-1, n-\ell) &= \sum_{\ell=0}^d \sum_{n=d}^N \frac{\binom{n-\ell}{d-\ell} \binom{N-n+\ell}{\ell}}{\binom{N}{d}} \pi(x-1, n-\ell), \\ &= \sum_{\ell=0}^d \sum_{n=d-\ell}^{N-\ell} \frac{\binom{n}{d-\ell} \binom{N-n}{\ell}}{\binom{N}{d}} \pi(x-1, n), \\ &= \sum_{n=0}^N \left[ \sum_{\ell=(d-n)^+}^{d \wedge (N-n)} \frac{\binom{n}{d-\ell} \binom{N-n}{\ell}}{\binom{N}{d}} \right] \pi(x-1, n), \\ &= \sum_{n=0}^N \pi(x-1, n), \end{aligned}$$

en reconnaissant la fonction de masse d'une loi hypergéométrique de paramètres  $N, N-n, d$ . □

**Remarque :** pour  $x = 1$ , il est facile de voir que l'on a en fait une égalité :

$$\sum_{n=0}^N \pi(1, n) = \pi(1, d) = \frac{N\rho}{d} \sum_{n=0}^N \pi(0, n) = \frac{N\rho}{d} \pi(0, 0).$$

**Corollaire 1.** *Pour tout  $x \geq 0$ , on a*

$$\sum_{n=0}^N \pi(x, n) \leq \left( \frac{N\rho}{d} \right)^x \pi(0, 0).$$

*Il vient en particulier*

$$\pi(0, 0) \geq \frac{1}{2}, \quad \sum_{x=2}^{\infty} \sum_{n=0}^N \pi(x, n) = \underset{\frac{N\rho}{d} \rightarrow 0}{O} \left( \left( \frac{N\rho}{d} \right)^2 \right) \quad \text{et} \quad \sum_{x=3}^{\infty} x \sum_{n=0}^N \pi(x, n) = \underset{\frac{N\rho}{d} \rightarrow 0}{O} \left( \left( \frac{N\rho}{d} \right)^3 \right).$$

À la racine du parallélisme

*Démonstration.* Le premier point découle du Lemme 1 par récurrence. Concernant le deuxième point, avec l'hypothèse  $N\rho/d \leq 1/2$ , on a

$$1 = \pi(0,0) + \sum_{x=1}^{+\infty} \sum_{n=0}^N \pi(x,n) \leq \pi(0,0) \left(1 + \sum_{x=1}^{+\infty} \left(\frac{N\rho}{d}\right)^x\right) = \pi(0,0) \frac{1}{1 - \frac{N\rho}{d}} \leq 2\pi(0,0)$$

Pour le reste, on a d'une part

$$\sum_{x=2}^{+\infty} \sum_{n=0}^N \pi(x,n) \leq \sum_{x=2}^{+\infty} \left(\frac{N\rho}{d}\right)^x = \frac{\left(\frac{N\rho}{d}\right)^2}{1 - \frac{N\rho}{d}} \leq 2 \left(\frac{N\rho}{d}\right)^2,$$

et d'autre part, en utilisant l'égalité  $\sum_{x=k}^{+\infty} x a^x = \frac{ka^k - (k-1)a^{k+1}}{(1-a)^2}$  pour  $|a| < 1$ ,

$$\sum_{x=3}^{+\infty} x \sum_{n=0}^N \pi(x,n) \leq \sum_{x=3}^{+\infty} x \left(\frac{N\rho}{d}\right)^x = \frac{3\left(\frac{N\rho}{d}\right)^3 - 2\left(\frac{N\rho}{d}\right)^4}{\left(1 - \frac{N\rho}{d}\right)^2} < 12 \left(\frac{N\rho}{d}\right)^3.$$

□

Le corollaire 1 peut s'interpréter comme suit : sachant que le taux de service du cluster vaut au minimum  $d\mu$ , le système traite les requêtes au moins aussi vite qu'une file  $M/M/1$  d'intensité d'arrivée  $N\lambda$  et de taux de service  $d\mu$ . En particulier, la probabilité d'avoir  $x$  requêtes décroît au moins comme une suite géométrique de raison  $\frac{N\lambda}{d\mu} = \frac{N\rho}{d}$ .

Nous pouvons maintenant revenir au calcul de  $\alpha_d(\rho)$ . D'après (4), on a,

$$\alpha_d(\rho) = \frac{1 - \rho}{\rho} \frac{d\mathbb{E}(\mathbf{X})}{N} = \frac{1 - \rho}{\frac{N\rho}{d}} \times \frac{\sum_{x=0}^{+\infty} x \sum_{n=0}^N \pi(x,n)}{\sum_{x=0}^{+\infty} \sum_{n=0}^N \pi(x,n)}.$$

En appliquant le Corollaire 1 puis (1), il vient

$$\begin{aligned} \alpha_d(\rho) &= \frac{1 - \rho}{\frac{N\rho}{d}} \times \frac{0 \times \pi(0,0) + 1 \times \pi(1,d) + 2 \times \sum_{n=0}^N \pi(2,n) + O\left(\left(\frac{N\rho}{d}\right)^3\right)}{\pi(0,0) + \pi(1,d) + O\left(\left(\frac{N\rho}{d}\right)^2\right)}, \\ &= \frac{1 - \rho}{\frac{N\rho}{d}} \times \frac{\frac{N\rho}{d} + \frac{2}{\pi(0,0)} \sum_{k=0}^d \pi(2,2d-k) + O\left(\left(\frac{N\rho}{d}\right)^3\right)}{1 + \frac{N\rho}{d} + O\left(\left(\frac{N\rho}{d}\right)^2\right)}. \end{aligned}$$

Pour chaque  $k = (2d - N)^+, \dots, d$ , on a d'après (1),

$$\pi(2,2d-k) = \frac{N\rho}{2d-k} \sum_{\ell=0}^d \frac{\binom{2d-k-\ell}{d-\ell} \binom{N-2d+k+\ell}{\ell}}{\binom{N}{d}} \pi(1,2d-k-\ell),$$

et puisque  $\pi(1,n) = 0$  si  $n \neq d$ , il vient

$$\pi(2,2d-k) = \frac{N\rho}{2d-k} \frac{\binom{d}{k} \binom{N-d}{d-k}}{\binom{N}{d}} \pi(1,d) = \frac{N\rho}{2d-k} \frac{\binom{d}{k} \binom{N-d}{d-k}}{\binom{N}{d}} \frac{N\rho}{d} \pi(0,0) = \frac{N\rho}{2d-k} p_k \frac{N\rho}{d} \pi(0,0),$$



où  $(p_k)_{k=0,\dots,d}$  dénote la fonction de masse de la loi hypergéométrique de paramètres  $N, d, d$  :  $p_k = 0$  si  $k < (2d - N)^+$ , et pour chaque  $k = (2d - N)^+, \dots, d$ ,

$$p_k = \frac{\binom{d}{k} \binom{N-d}{d-k}}{\binom{N}{d}}.$$

En injectant ceci dans l'expression de  $\alpha_d(\rho)$ , on obtient ainsi

$$\alpha_d(\rho) = (1 - \rho) \frac{1 + \frac{N\rho}{d} \sum_{k=0}^d \frac{1}{1 - \frac{k}{2d}} p_k + \mathcal{O}\left(\left(\frac{N\rho}{d}\right)^2\right)}{1 + \frac{N\rho}{d} + \mathcal{O}\left(\left(\frac{N\rho}{d}\right)^2\right)}. \quad (7)$$

Le reste de la preuve consiste à approximer la somme faisant intervenir  $(p_k)_{k=0,\dots,d}$  dans (7), dans trois régimes différents.

### B.1 Cas $d$ petit devant $\sqrt{N}$

On a en particulier  $d \leq N/2$ , de sorte que  $2d - N \leq 0$ . Le Lemme suivant permet de comparer  $(p_k)_{k=0,\dots,d}$  avec la fonction de masse d'une loi de Poisson de paramètre  $d^2/N$ .

**Lemme 2.** Pour chaque  $k = 0, \dots, d$ , on a

$$p_k = \frac{1}{k!} \times \prod_{\ell=0}^{k-1} \frac{(d-\ell)^2}{N-\ell} \times \prod_{\ell=k}^{d-1} \left(1 - \frac{d-k}{N-\ell}\right).$$

*Démonstration.* Soit  $k = 0, \dots, d$ . On a d'une part

$$\binom{d}{k} = \frac{d!}{k!(d-k)!} = \frac{1}{k!} \times \prod_{\ell=0}^{k-1} (d-\ell),$$

et d'autre part

$$\begin{aligned} \binom{N-d}{d-k} &= \frac{(N-d)!}{(d-k)!(N-2d+k)!}, \\ &= \frac{N!}{d!(N-d)!} \times \frac{\frac{d!}{(d-k)!}}{\frac{N!}{(N-k)!}} \times \frac{\frac{(N-d)!}{(N-d-(d-k))!}}{\frac{(N-k)!}{(N-k-(d-k))!}}, \\ &= \binom{N}{d} \times \prod_{\ell=0}^{k-1} \left(\frac{d-\ell}{N-\ell}\right) \times \prod_{\ell=0}^{d-k-1} \left(\frac{N-d-\ell}{N-k-\ell}\right), \\ &= \binom{N}{d} \times \prod_{\ell=0}^{k-1} \left(\frac{d-\ell}{N-\ell}\right) \times \prod_{\ell=0}^{d-k-1} \left(1 - \frac{d-k}{N-k-\ell}\right), \\ &= \binom{N}{d} \times \prod_{\ell=0}^{k-1} \left(\frac{d-\ell}{N-\ell}\right) \times \prod_{\ell=k}^{d-1} \left(1 - \frac{d-k}{N-\ell}\right). \end{aligned}$$

On obtient le résultat attendu en combinant ces deux égalités dans l'expression de  $p_k$ . □

En écrivant  $p_0 = 1 - p_1 - \sum_{k=2}^d p_k$ , il vient

$$\begin{aligned} \sum_{k=0}^d \frac{1}{1 - \frac{k}{2d}} p_k &= 1 \times \left(1 - p_1 - \sum_{k=2}^d p_k\right) + \frac{1}{1 - \frac{1}{2d}} p_1 + \sum_{k=2}^d \frac{1}{1 - \frac{k}{2d}} p_k, \\ &= 1 + \left(\frac{1}{1 - \frac{1}{2d}} - 1\right) p_1 + \sum_{k=2}^d \left(\frac{1}{1 - \frac{k}{2d}} - 1\right) p_k. \end{aligned}$$

À la racine du parallélisme

D'après le Lemme 2, on peut réécrire le second terme sous la forme

$$\left(\frac{1}{1-\frac{1}{2d}} - 1\right) p_1 = \frac{1}{2d-1} \frac{d^2}{N} \prod_{\ell=0}^{d-1} \left(1 - \frac{d-1}{N-\ell}\right) = \frac{1}{2-\frac{1}{d}} \frac{d}{N} \prod_{\ell=0}^{d-1} \left(1 - \frac{d-1}{N-\ell}\right).$$

De plus, la somme pour  $k = 2, \dots, d$  s'écrit

$$\sum_{k=2}^d \left(\frac{1}{1-\frac{k}{2d}} - 1\right) p_k = \sum_{k=2}^d \frac{k}{2d-k} p_k \leq \frac{d}{2d-d} \sum_{k=2}^d p_k = \sum_{k=2}^d p_k,$$

et d'après le Lemme 2,

$$\sum_{k=2}^d p_k \leq \sum_{k=2}^d \frac{1}{k!} \left(\frac{d^2}{N}\right)^k \leq e^{\frac{d^2}{N}} - \frac{d^2}{N} - 1,$$

ce qui donne finalement

$$\sum_{k=2}^d \left(\frac{1}{1-\frac{k}{2d}} - 1\right) p_k = O\left(\frac{d^4}{N^2}\right).$$

Une manière d'interpréter ce résultat est que quand  $d$  est petit devant  $\sqrt{N}$ , les états avec 2 collisions ou plus sont très rares et peuvent être négligés. Le numérateur dans (7) s'écrit donc

$$\begin{aligned} & 1 + \frac{N\rho}{d} \left(1 + \frac{1}{2-\frac{1}{d}} \frac{d}{N} \prod_{\ell=0}^{d-1} \left(1 - \frac{d-1}{N-\ell}\right) + O\left(\frac{d^4}{N^2}\right)\right) + O\left(\left(\frac{N\rho}{d}\right)^2\right) \\ &= 1 + \frac{N\rho}{d} + \rho \frac{1}{2-\frac{1}{d}} \prod_{\ell=0}^{d-1} \left(1 - \frac{d-1}{N-\ell}\right) + O\left(\frac{d^3\rho}{N} + \left(\frac{N\rho}{d}\right)^2\right). \end{aligned}$$

En injectant ceci dans (7), on obtient

$$\alpha_d(\rho) = (1-\rho) \frac{1 + \frac{N\rho}{d} + \rho \frac{1}{2-\frac{1}{d}} \prod_{\ell=0}^{d-1} \left(1 - \frac{d-1}{N-\ell}\right) + O\left(\frac{d^3\rho}{N}\right) + O\left(\left(\frac{N\rho}{d}\right)^2\right)}{1 + \frac{N\rho}{d} + O\left(\left(\frac{N\rho}{d}\right)^2\right)}.$$

Le coefficient qui suit  $\rho/(2-1/d)$  est proche de 1 lorsque  $d$  est petit devant  $N$ . On écrit alors l'approximation suivante, valable pour  $\frac{N\rho}{d}$  petit et  $d$  petit devant  $\sqrt{N}$ .

$$\alpha_d(\rho) \simeq (1-\rho) \left(1 + \frac{\rho}{2-\frac{1}{d}}\right) \simeq \frac{1-\rho}{1-\frac{\rho}{2-\left(\frac{1}{d}+\frac{1}{N}\right)}}.$$

## B.2 Cas $d$ proche de $\sqrt{N}$

On utilise ici le développement limité de  $1/(1-x)$  quand  $x$  tend vers 0 pour montrer que l'on peut négliger les termes correspondant à  $k \geq 3$  dans la somme. Plus précisément, pour tout  $x \in [0, \frac{1}{2}]$ , on a

$$\left| \frac{1}{1-x} - (1+x+x^2) \right| \leq 2x^3,$$

ce qui implique

$$\left| \sum_{k=0}^d \frac{1}{1-\frac{k}{2d}} p_k - \sum_{k=0}^d \left(1 + \frac{k}{2d} + \frac{k^2}{4d^2}\right) p_k \right| \leq \sum_{k=0}^d \left| \frac{1}{1-\frac{k}{2d}} - \left(1 + \frac{k}{2d} + \frac{k^2}{4d^2}\right) \right| p_k \leq 2 \sum_{k=0}^d \frac{k^3}{8d^3} p_k.$$

Il vient ainsi

$$\begin{aligned} \sum_{k=0}^d \frac{1}{1-\frac{k}{2d}} p_k &= \sum_{k=0}^d \left(1 + \frac{k}{2d} + \frac{k^2}{4d^2}\right) p_k + O\left(\sum_{k=0}^d \frac{k^3}{d^3} p_k\right), \\ &= 1 + \frac{1}{2d} \mathbb{E}(\mathbf{K}) + \frac{1}{4d^2} \mathbb{E}(\mathbf{K}^2) + O\left(\frac{1}{d^3} \mathbb{E}(\mathbf{K}^3)\right). \end{aligned}$$

où  $\mathbf{K}$  dénote une variable aléatoire de distribution  $p_k$ . On sait par ailleurs que

$$\mathbb{E}(\mathbf{K}) = \frac{d^2}{N}, \quad \mathbb{E}(\mathbf{K}^2) = \frac{d^4}{N^2} + \frac{d^2(N-d)^2}{N^2(N-1)},$$

et en notant  $\mathbf{L}$  une variable aléatoire de distribution hypergéométrique de paramètres  $N-1, d-1, d-1$ ,

$$\begin{aligned} \mathbb{E}(\mathbf{K}^3) &= \frac{d^2}{N} \mathbb{E}((1+\mathbf{L})^2) = \frac{d^2}{N} (1+2\mathbb{E}(\mathbf{L}) + \mathbb{E}(\mathbf{L}^2)), \\ &= \frac{d^2}{N} \left(1+2\frac{(d-1)^2}{N-1} + \frac{(d-1)^4}{(N-1)^2} + \frac{(d-1)^2(N-d)^2}{(N-1)^2(N-2)}\right). \end{aligned}$$

Il vient ainsi

$$\begin{aligned} \sum_{k=0}^d \frac{1}{1-\frac{k}{2d}} p_k &= 1 + \frac{1}{2d} \frac{d^2}{N} + \frac{1}{4d^2} \left(\frac{d^4}{N^2} + \frac{d^2(N-d)^2}{N^2(N-1)}\right) \\ &\quad + O\left(\frac{1}{d^3} \frac{d^2}{N} \left(1+2\frac{(d-1)^2}{N-1} + \frac{(d-1)^4}{(N-1)^2} + \frac{(d-1)^2(N-d)^2}{(N-1)^2(N-2)}\right)\right), \end{aligned}$$

d'où

$$\sum_{k=0}^d \frac{1}{1-\frac{k}{2d}} p_k = 1 + \frac{1}{2} \frac{d}{N} + \frac{1}{4} \frac{d^2}{N^2} + \frac{1}{4} \frac{(N-d)^2}{N^2(N-1)} + O\left(\frac{1}{dN}\right).$$

Le numérateur de (7) s'écrit ainsi

$$\begin{aligned} &1 + \frac{N\rho}{d} \sum_{k=0}^d \frac{1}{1-\frac{k}{2d}} p_k + O\left(\left(\frac{N\rho}{d}\right)^2\right) \\ &= 1 + \frac{N\rho}{d} \left(1 + \frac{1}{2} \frac{d}{N} + \frac{1}{4} \frac{d^2}{N^2} + \frac{1}{4} \frac{(N-d)^2}{N^2(N-1)}\right) + O\left(\frac{\rho}{d^2} + \left(\frac{N\rho}{d}\right)^2\right), \\ &= 1 + \frac{N\rho}{d} + \frac{1}{2} \rho \left(1 + \frac{1}{2} \left(\frac{d}{N} + \frac{1}{d} \frac{(N-d)^2}{N(N-1)}\right)\right) + O\left(\frac{\rho}{d^2} + \left(\frac{N\rho}{d}\right)^2\right). \end{aligned}$$

On obtient donc pour  $\alpha_d(\rho)$  :

$$\alpha_d(\rho) = (1-\rho) \frac{1 + \frac{N\rho}{d} + \frac{1}{2} \rho \left(1 + \frac{1}{2} \left(\frac{d}{N} + \frac{1}{d} \frac{(N-d)^2}{N(N-1)}\right)\right) + O\left(\frac{\rho}{d^2} + \left(\frac{N\rho}{d}\right)^2\right)}{1 + \frac{N\rho}{d} + O\left(\left(\frac{N\rho}{d}\right)^2\right)}.$$

En faisant tendre  $\frac{N\rho}{d}$  vers 0, on obtient à nouveau l'équation (5).

À la racine du parallélisme

### B.3 Cas $d$ proche de $N$

On suppose en particulier  $d \geq N/2$ , de sorte que  $2d - N \geq 0$ . On utilise le Lemme suivant pour se ramener au cas où  $d$  est petit devant  $N$ . Il signifie que le nombre de collisions en plus du nombre minimum  $2d - N$  est distribué selon une loi hypergéométrique de paramètres  $N, N - d, N - d$ .

**Lemme 3.** Pour chaque  $k = 2d - N, \dots, d$ , on a

$$\frac{\binom{d}{k} \binom{N-d}{d-k}}{\binom{N}{d}} = \frac{\binom{N-d}{\ell} \binom{N-(N-d)}{(N-d)-\ell}}{\binom{N}{N-d}}$$

où  $\ell = k - (2d - N)$ .

*Démonstration.* On a directement  $\binom{N}{d} = \binom{N}{N-d}$ . Soient  $k = 2d - N, \dots, d$  et  $\ell = k - (2d - N)$ . On a

$$\binom{d}{k} \binom{N-d}{d-k} = \binom{d}{\ell + 2d - N} \binom{N-d}{d - (\ell + 2d - N)} = \binom{d}{\ell + 2d - N} \binom{N-d}{N-d-\ell},$$

puis par symétrie des coefficients binomiaux :

$$\binom{d}{k} \binom{N-d}{d-k} = \binom{d}{N-d-\ell} \binom{N-d}{\ell} = \binom{N-(N-d)}{(N-d)-\ell} \binom{N-d}{\ell}.$$

□

On peut donc appliquer la même méthode que dans l'Annexe B.1 en supposant  $N - d$  petit devant  $\sqrt{N}$ . Plus précisément, en notant

$$p'_\ell = \frac{\binom{N-d}{\ell} \binom{N-(N-d)}{(N-d)-\ell}}{\binom{N}{N-d}}$$

pour chaque  $\ell = 0, \dots, N - d$ , on obtient

$$\begin{aligned} \sum_{k=2d-N}^d \frac{1}{1 - \frac{k}{2d}} p_k &= \frac{1}{1 - \frac{2d-N}{2d}} \left( 1 - p_{2d-N+1} - \sum_{k=2d-N+2}^d p_k \right) + \frac{1}{1 - \frac{2d-N+1}{2d}} p_{2d-N+1} + \sum_{k=2d-N+2}^d \frac{1}{1 - \frac{k}{2d}} p_k, \\ &= \frac{1}{1 - \frac{2d-N}{2d}} + \left( \frac{1}{1 - \frac{2d-N+1}{2d}} - \frac{1}{1 - \frac{2d-N}{2d}} \right) p_{2d-N+1} + \sum_{k=2d-N+2}^d \left( \frac{1}{1 - \frac{k}{2d}} - \frac{1}{1 - \frac{2d-N}{2d}} \right) p_k, \\ &= \frac{1}{1 - \frac{1}{2} \frac{d}{N} + \frac{(N-d)^2}{2Nd}} + \frac{2d}{N(N-1)} p'_1 + \sum_{\ell=2}^d 2d \left( \frac{1}{N-\ell} - \frac{1}{N} \right) p'_\ell, \end{aligned}$$

en remarquant que  $2d - N = \frac{d^2}{N} - \frac{(N-d)^2}{N}$ . D'après le Lemme 2, on a

$$\frac{2d}{N(N-1)} p'_1 = \frac{2d}{N(N-1)} \frac{(N-d)^2}{N} \times \prod_{\ell=0}^{N-d-1} \left( 1 - \frac{N-d-1}{N-\ell} \right) \leq \frac{2d}{N(N-1)} \frac{(N-d)^2}{N},$$

et

$$\begin{aligned}
 \sum_{\ell=2}^d 2d \left( \frac{1}{N-\ell} - \frac{1}{N} \right) p'_\ell &\leq 2d \left( \frac{1}{N-d} - \frac{1}{N} \right) \sum_{\ell=2}^d p'_\ell = \frac{2d^2}{N(N-d)} \sum_{\ell=2}^d p'_\ell, \\
 &\leq \frac{2d^2}{N(N-d)} \sum_{\ell=2}^d \frac{1}{\ell!} \left( \frac{(N-d)^2}{N} \right)^\ell, \\
 &\leq \frac{2d^2}{N(N-d)} \left( e^{\frac{(N-d)^2}{N}} - \frac{(N-d)^2}{N} - 1 \right), \\
 &\leq \frac{2d^2}{N(N-d)} \frac{1}{2} \frac{(N-d)^4}{N^2} = \frac{d^2(N-d)^3}{N^3}, \\
 &\leq \frac{(N-d)^3}{N^2}.
 \end{aligned}$$

Il vient ainsi

$$\begin{aligned}
 \sum_{k=2d-N}^d \frac{1}{1 - \frac{k}{2d}} p_k &= \frac{1}{1 - \frac{1}{2} \frac{d}{N} + \frac{(N-d)^2}{2Nd}} + O\left(\frac{(N-d)^3}{N^2}\right), \\
 &= 1 + \frac{1}{2} \frac{d}{N} \frac{1 - \frac{(N-d)^2}{d^2}}{1 - \frac{1}{2} \frac{d}{N} + \frac{(N-d)^2}{2Nd}} + O\left(\frac{(N-d)^3}{N^2}\right), \\
 &= 1 + \frac{d}{N} \frac{1 - \frac{(N-d)^2}{d^2}}{2 - \frac{d}{N} + \frac{(N-d)^2}{Nd}} + O\left(\frac{(N-d)^3}{N^2}\right).
 \end{aligned}$$

Finalement, on obtient d'après (7)

$$\alpha_d(\rho) = (1-\rho) \frac{1 + \frac{N\rho}{d} + \rho \frac{1 - \frac{(N-d)^2}{d^2}}{2 - \frac{d}{N} + \frac{(N-d)^2}{Nd}} + O\left(\rho \frac{(N-d)^3}{Nd} + \left(\frac{N\rho}{d}\right)^2\right)}{1 + \frac{N\rho}{d} + O\left(\left(\frac{N\rho}{d}\right)^2\right)}.$$

En faisant tendre  $\frac{N\rho}{d}$  vers 0, on obtient à nouveau l'équation (5).

## C Preuve de la Proposition 2 de la Partie 5

Soit  $2 \leq d \leq N-1$ . D'après (4), on a

$$\alpha_d(\rho) = \frac{d}{N} \sum_{n=d}^N \frac{1-\rho}{\binom{N-1}{d-1} - \rho} = \frac{d}{N} + \frac{d}{N} (1-\rho) \sum_{n=d}^{N-1} \frac{1}{\binom{N-1}{d-1} - \rho}.$$

La somme peut être réécrite

$$\sum_{n=d}^{N-1} \frac{1}{\binom{N-1}{d-1} - \rho} = \sum_{n=d}^{N-1} \frac{\binom{n-1}{d-1}}{\binom{N-1}{d-1} - \binom{n-1}{d-1}} \frac{1}{1 + \frac{\binom{n-1}{d-1}}{\binom{N-1}{d-1} - \binom{n-1}{d-1}} (1-\rho)}.$$

En utilisant la majoration suivante pour chaque  $n = d, \dots, N-1$ , on a

$$\frac{\binom{n-1}{d-1}}{\binom{N-1}{d-1} - \binom{n-1}{d-1}} \leq \frac{\binom{N-2}{d-1}}{\binom{N-1}{d-1} - \binom{N-2}{d-1}} = \frac{N-d}{d-1},$$

À la racine du parallélisme

on obtient après un développement limité :

$$\begin{aligned} \sum_{n=d}^{N-1} \frac{1}{\frac{\binom{N-1}{d-1}}{\binom{n-1}{d-1}} - \rho} &= \sum_{n=d}^{N-1} \frac{\binom{n-1}{d-1}}{\binom{N-1}{d-1} - \binom{n-1}{d-1}} \left( 1 + O\left( (1-\rho) \frac{N-d}{d} \right) \right), \\ &= O\left( \frac{(N-d)^2}{d} + (1-\rho) \frac{(N-d)^3}{d^2} \right). \end{aligned}$$

Finalement,

$$\alpha_d(\rho) = \frac{d}{N} + O\left( (1-\rho) \frac{(N-d)^2}{N} + (1-\rho)^2 \frac{(N-d)^3}{Nd} \right) \simeq \frac{d}{N}.$$