



HAL
open science

The critical Karp-Sipser core of random graphs

Thomas Budzinski, Alice Contat, Nicolas Curien

► **To cite this version:**

Thomas Budzinski, Alice Contat, Nicolas Curien. The critical Karp-Sipser core of random graphs. 2022. hal-04285051

HAL Id: hal-04285051

<https://hal.science/hal-04285051>

Preprint submitted on 14 Nov 2023

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

THE CRITICAL KARP–SIPSER CORE OF RANDOM GRAPHS

Thomas BUDZINSKI* & Alice CONTAT† & Nicolas CURIEN‡

Abstract

We study the Karp–Sipser core of a random graph made of a configuration model with vertices of degree 1, 2 and 3. This core is obtained by recursively removing the leaves as well as their unique neighbors in the graph. We settle a conjecture of Bauer & Golinelli [2] and prove that at criticality, the Karp–Sipser core has size $\approx \text{Cst} \cdot \vartheta^{-2} \cdot n^{3/5}$ where ϑ is the hitting time of the curve $t \mapsto \frac{1}{t^2}$ by a linear Brownian motion started at 0. Our proof relies on a detailed multi-scale analysis of the Markov chain associated to the Karp–Sipser leaf-removal algorithm close to its extinction time.

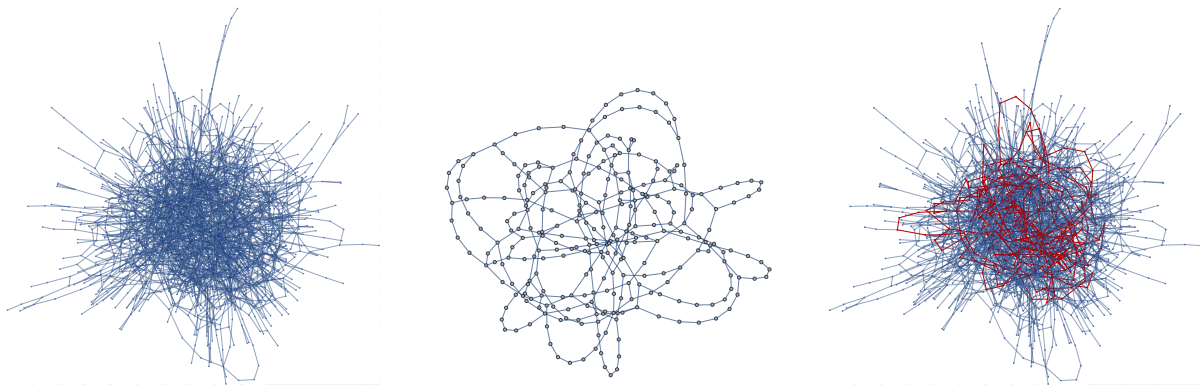


Figure 1: **(Left)**. The giant component of an Erdős–Rényi random graph $G(n, \frac{c}{n})$ with $n = 2000$ on the left and **(Middle)** its Karp–Sipser core. **(Right)**. The Karp–Sipser core in red inside the original graph.

1 Introduction

The Karp–Sipser algorithm. Let \mathfrak{g} be a finite graph. The Karp–Sipser algorithm [11] consists in removing recursively the vertices of degree 1 in \mathfrak{g} as well as their unique neighbors, see Figure 2. The initial motivation of Karp & Sipser for considering this algorithm is that the leaves¹ and isolated vertices removed during this process form an independent set of \mathfrak{g} which has very high density. We recall that an independent set in \mathfrak{g} is a subset of vertices, no two of which are adjacent. The problem

*ENS de Lyon and CNRS.

†Université Paris-Saclay.

‡Université Paris-Saclay.

¹Here and in the rest of the paper, the concept of leaf is a dynamical concept, as a vertex in the initial graph which is not a leaf may become one later.

thomas.budzinski@ens-lyon.fr
alice.contat@universite-paris-saclay.fr
nicolas.curien@gmail.com

of finding an independent set of maximal size is in general a NP-hard problem, and the Karp–Sipser algorithm provides a fair lower bound (it is furthermore “optimal” as long as there are remaining leaves in the graph).

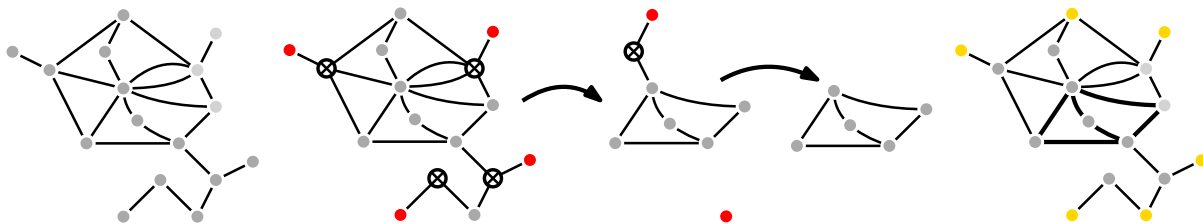


Figure 2: Illustration of the Karp–Sipser algorithm. The first 4 figures show the initial graph, as well as the recursive deletion process of the leaves (in red) together with their unique neighbor (crosses), until no leaf is left: we then obtain the Karp–Sipser core (fourth figure). On the right, the initial graph is represented together with the Karp–Sipser core in thick lines and the independent set formed by the removed “leaves” in yellow.

The Karp–Sipser core of random graphs. A striking property of the leaf-removal process is its Abelian property: whatever the order in which we decide to recursively remove the leaves and their neighbors, we always obtain the same subgraph of \mathbf{g} (with no leaves) which we will call the *Karp–Sipser core* of \mathbf{g} and denote by $\text{KSCore}(\mathbf{g})$, see [2, Appendix] or [12, Section 1.6.1]. Beware that the above notion differs from the usual k -core of a graph², see Section 5. By the above remark, the Karp–Sipser algorithm creates an independent set (the leaves removed during the algorithm) whose size is within at most $|\text{KSCore}(\mathbf{g})|$ from the maximal size of an independent set in \mathbf{g} .

The performance of the Karp–Sipser algorithm on the Erdős–Rényi random graph $G(n, \frac{c}{n})$ has been analyzed in the pioneer work [11] and later refined in the breakthrough [1] which established a phase transition as $n \rightarrow \infty$ depending on the value of c :

- if $c < e$, then as $n \rightarrow \infty$, the size $|\text{KSCore}(G(n, \frac{c}{n}))|$ is of order $O(1)$;
- if $c > e$, then as $n \rightarrow \infty$, the size $|\text{KSCore}(G(n, \frac{c}{n}))|$ is of order n .

Those works have later been extended to the configuration model [4, 10]. However, the careful analysis of the critical case $c = e$ was open as of today to the best of our knowledge. In [2], based on numerical simulations, the physicists Bauer & Golinelli predicted that $|\text{KSCore}(G(n, \frac{e}{n}))|$ should be of order $n^{3/5}$. The main result of this work (Theorem 2) is to settle this conjecture in the case of a random graph with degrees 1, 2 and 3.

Model and results. In this paper we shall consider a random graph model closely related to $G(n, \frac{c}{n})$ but for which the analysis of the Karp–Sipser algorithm is simpler. Namely, we fix a sequence of numbers $\mathbf{d}^n = (d_1^n, d_2^n, d_3^n)_{n \geq 1}$ such that

$$n = d_1^n + 2d_2^n + 3d_3^n \text{ is even.}$$

²The k -core of \mathbf{g} is the largest subset V of its vertices such that for any $v \in V$, the induced degree of v within V is at least k .

We imagine \mathbf{d}^n as the number of vertices of degree 1, 2 and 3 and consider a random multi-graph $\text{CM}(\mathbf{d}^n)$ sampled by pairing the edges emanating from the $d_1^n + d_2^n + d_3^n$ vertices uniformly at random. This is a special instance of the so-called configuration model introduced by Bollobas [5], see [18] for background. In the rest of the paper we shall further assume that

$$\frac{d_1^n}{n} \xrightarrow{n \rightarrow \infty} p_1, \quad \frac{2d_2^n}{n} \xrightarrow{n \rightarrow \infty} p_2, \quad \text{and} \quad \frac{3d_3^n}{n} \xrightarrow{n \rightarrow \infty} p_3, \quad (1)$$

so that the proportion of half-edges which are incident to a vertex of degree i is p_i . Our goal will be to analyze $\text{KSCore}(\text{CM}(\mathbf{d}^n))$. A phase transition has been observed in [10] for the size of the Karp–Sipser core but its location depending on (p_1, p_2, p_3) was not explicit. Our first contribution is to make this threshold precise. For a graph \mathbf{g} , we will write $|\mathbf{g}|$ for twice the number of edges of \mathbf{g} , and call this quantity the *size* of \mathbf{g} . If (u_n) is a sequence of positive numbers and (X_n) a sequence of random variables, we will write $X_n = O_{\mathbb{P}}(u_n)$ if $(u_n^{-1}X_n)$ is tight, and we will write $X_n = o_{\mathbb{P}}(u_n)$ if $u_n^{-1}X_n$ converges to 0 in probability.

Theorem 1 (Explicit phase transition). *Under the assumptions (1), let*

$$\Theta = (p_3 - p_1)^2 - 4p_1. \quad (2)$$

- **Subcritical phase.** *If $\Theta < 0$, then as $n \rightarrow \infty$ we have*

$$|\text{KSCore}(\text{CM}(\mathbf{d}^n))| = O_{\mathbb{P}}(\log^2 n).$$

- **Supercritical phase.** *If $\Theta > 0$, then*

$$n^{-1} \cdot |\text{KSCore}(\text{CM}(\mathbf{d}^n))| \xrightarrow[n \rightarrow \infty]{(\mathbb{P})} \frac{4\Theta}{3 + \Theta}.$$

- **Critical phase.** *If $\Theta = 0$, then $|\text{KSCore}(\text{CM}(\mathbf{d}^n))| = o_{\mathbb{P}}(n)$.*

Sketch of proof of the phase transition. The proof of this theorem uses classical techniques. We shall reveal the random graph $\text{CM}(\mathbf{d}^n)$ by pairing its half-edges two-by-two as we perform the Karp–Sipser leaf removal algorithm (a.k.a. peeling algorithm). More precisely, when we remove a leaf, we reveal its neighbor in the graph and remove it as well, which decreases the degrees of some other vertices. During this process, the number of remaining vertices of degree 1, 2 and 3 evolves as an $(\mathbb{Z}_{\geq 0})^3$ -valued Markov chain with explicit probability transitions. This is, of course, a recurrent idea in random graph theory and has already been used many times for the Karp–Sipser algorithm itself [11, 1]. More precisely, we shall erase leaves uniformly at random one-by-one (in contrast with [10], where all possible leaves are erased at each round) and use the *differential equation method* [19] to prove that the renormalized number of vertices of degree 1, 2 and 3 is well approximated by a differential equation on \mathbb{R}^3 for which we are able to find explicit solutions. In a sense, this returns to the roots of this method since it was Karp & Sipser [11] who first introduced it in the context of random graphs following earlier works of Kurtz [13] in population models.

Remark (A spectral parallel to the Karp–Sipser phase transition). The *nullity* of a graph is the multiplicity of 0 in the spectrum of its adjacency matrix. It is easy to see that the leaf-removal process on a graph \mathbf{g} leaves its *nullity* invariant and so the Karp–Sipser algorithm can also be used to study the later, see [3, 17]. The phase transition for the emergence of a Karp–Sipser core of

positive proportion in $G(n, \frac{\varepsilon}{n})$ has a parallel phase transition³ for the emergence of extended states (an absolutely continuous part) at zero in $G(n, \frac{\varepsilon}{n})$, see [3, 6]. We wonder whether a similar result holds true for the configuration models we study.

We now turn to the detailed analysis of the *critical case* which is the main goal of our work. For this we fix a particular degree sequence $\mathbf{d}_{\text{crit}}^n = (d_{1,c}^n, d_{2,c}^n, d_{3,c}^n)$ such that $d_{1,c}^n + 3d_{3,c}^n = n$ is even (to be able to perform the configuration model) and

$$d_{1,c}^n = n(1 - \frac{\sqrt{3}}{2}) + O(1), \quad 2d_{2,c}^n = 0, \quad \text{and} \quad 3d_{3,c}^n = n\frac{\sqrt{3}}{2} + O(1). \quad (3)$$

In particular we have $\Theta = (\sqrt{3} - 1)^2 - 4(1 - \frac{\sqrt{3}}{2}) = 0$ so we are indeed in the critical case of Theorem 1. By definition, the core $\text{KSCore}(\text{CM}(\mathbf{d}_{\text{crit}}^n))$ has only vertices of degrees 2 or 3. Our main result is then the following:

Theorem 2 (Geometry of the critical Karp–Sipser core). *Let $D_2(n)$ (resp. $D_3(n)$) be the total number of half-edges attached to a vertex of degree 2 (resp. 3) in $\text{KSCore}(\text{CM}(\mathbf{d}_{\text{crit}}^n))$. Then we have*

$$\begin{pmatrix} n^{-3/5} \cdot D_2(n) \\ n^{-2/5} \cdot D_3(n) \end{pmatrix} \xrightarrow[n \rightarrow \infty]{(d)} \begin{pmatrix} 3^{-3/5} 2^{14/5} \cdot \vartheta^{-2} \\ 3^{-2/5} 2^{16/5} \cdot \vartheta^{-3} \end{pmatrix},$$

where $\vartheta = \inf\{t \geq 0 : B_t = t^{-2}\}$, for a standard linear Brownian motion $(B_t : t \geq 0)$ started from 0. Moreover, conditionally on $(D_2(n), D_3(n))$, the graph $\text{KSCore}(\text{CM}(\mathbf{d}_{\text{crit}}^n))$ is a configuration model.

Remark (Bauer & Golinelli’s prediction). The above theorem confirms a long-standing prediction of Bauer & Golinelli [2] stated in the case of the Erdős–Rényi random graph: based on Monte-Carlo simulations they proposed a few possible sets of critical exponents [2, Table 1] and our theorem confirms their prediction. See also [8, 12] for later developments.

Note that our assumptions on the initial degree sequence are much stronger than for Theorem 1 since the size of the critical core is quite sensitive to initial conditions. Our proof still works if the error $O(1)$ is replaced by $O(n^{1/2})$, and the result should remain true as long as the initial error is $o(n^{3/5})$, see Section 5 for a discussion on the near-critical regime. Although our main result only considers the graph $\text{CM}(\mathbf{d}_{\text{crit}}^n)$, we believe that the above limiting result holds for a large variety of random graphs which are critical for the Karp–Sipser algorithm. In particular, we expect a similar result for configuration models with bounded degrees and for the Erdős–Rényi graph $G(n, \frac{\varepsilon}{n})$, but the number of vertices of degree $2 \leq k \leq 5$ in the core should be of order $n^{(5-k)/5}$. In particular, we conjecture that there are no vertices of degree 6 or more in $\text{KSCore}(G(n, \frac{\varepsilon}{n}))$.

Ideas of proof. The proof of Theorem 2 uses the same Markov chain as the one used to study the phase transition. The difference is that we need to study the behaviour of this chain right before its extinction, at a scale much finer than n . More precisely, we can expect from the differential equation approximation that εn steps before extinction, the number of vertices of unmatched degrees 1, 2 and 3 are respectively of order $\varepsilon^2 n$, εn and $\varepsilon^{3/2} n$. On the other hand, a variance computation shows that the fluctuations of the number of vertices of degree 1 are of order $\varepsilon^{3/4} \sqrt{n}$. Finally, the time at which we can expect the Markov chain to terminate is the time where the fluctuations exceed the expected

³Unfortunately this does not seem to be an easy corollary of the “geometric” phase transition.

value, that is at $\varepsilon = n^{-2/5}$. However, checking that the differential equation approximation remains good until that scale requires some careful control of the Markov chain accross scales. In particular, the reason why the fluctuations become much smaller than \sqrt{n} in the end of the process is that the drift of our Markov chain induces a “self-correcting” effect.

Acknowledgments. The last two authors were supported by ERC 740943 GeoBrown and by ANR RanTanPlan. The first author is grateful to the Laboratoire de Mathématiques d’Orsay, where most of this work was done, for its hospitality. We warmly thank Matthieu Jonckheere for a stimulation discussion about [10] and Justin Salez for enlightening explanations about maximal matchings and independent sets in random graphs.

2 Karp–Sipser exploration of the configuration model

As we mentioned in the introduction the main idea (already present in [11, 1, 10, 4, 12]) is to explore the random configuration model $\text{CM}(\mathbf{d}^n)$ at the same time as we run the Karp–Sipser algorithm to discover its core. Let us explain this in details. Fix a degree sequence $\mathbf{d}^n = (d_1^n, d_2^n, d_3^n)$ such that $n = d_1^n + 2d_2^n + 3d_3^n$ is even. We shall expose the $\frac{n}{2}$ edges of $\text{CM}(\mathbf{d}^n)$ one by one and create a process

$$(X_k^n, Y_k^n, Z_k^n : k \geq 0)$$

where X^n, Y^n, Z^n represent respectively the number of unmatched half-edges linked to vertices of *unmatched degree*⁴ 1, 2, 3. The process of the sum is denoted by $S^n = X^n + Y^n + Z^n$. In particular, we always have $(X_0^n, Y_0^n, Z_0^n) = (d_1^n, 2d_2^n, 3d_3^n)$ and $S_0^n = n$ with our conventions.

As long as $X_k^n > 0$, the process evolves as follows. Since $X_k^n > 0$, there are still vertices of unmatched degree 1. We pick ℓ (for leaf) one of these vertices uniformly at random and reveal its neighbor v in the graph. Now, in the Karp–Sipser algorithm this vertex is “destroyed” so we shall erase v from the configuration *as well as the connections it has with other vertices of the graph*. More precisely, we reveal the neighbors of v in $\text{CM}(\mathbf{d}^n)$ and erase all the connections we create when doing so. In particular, if v is connected to a vertex $w \neq \ell$ of unmatched degree d via i edges, then after the operation w becomes a vertex of unmatched degree $d - i$. After that, the vertices of unmatched degree 0 are simply removed. We listed all 13 combinatorial possibilities (recall that our vertices have degree 1, 2 or 3) in Figure 3. The stopping time of the algorithm is

$$\theta^n := \inf\{k \geq 0 : X_k^n = 0\}.$$

Finally, we extend the process (X^n, Y^n, Z^n) to any k by setting $(X_k^n, Y_k^n, Z_k^n) = (X_{\theta^n}^n, Y_{\theta^n}^n, Z_{\theta^n}^n)$ for $k \geq \theta^n$. We denote by $(\mathcal{F}_k)_{k \geq 0}$ the natural filtration generated by this exploration. The starting point of our investigations is the following.

Proposition 1. *The process $(X_k^n, Y_k^n, Z_k^n)_{0 \leq k \leq \theta^n}$ is a Markov process whose probability transitions are described in Figure 3. Furthermore, for any stopping time τ , the remaining pairing of the unmatched edges conditionally on \mathcal{F}_τ is uniform.*

Proof. This is standard: the above exploration procedure of $\text{CM}(\mathbf{d}^n)$ is Markovian and preserves the fact that the remaining pairing of edges is uniform. \square

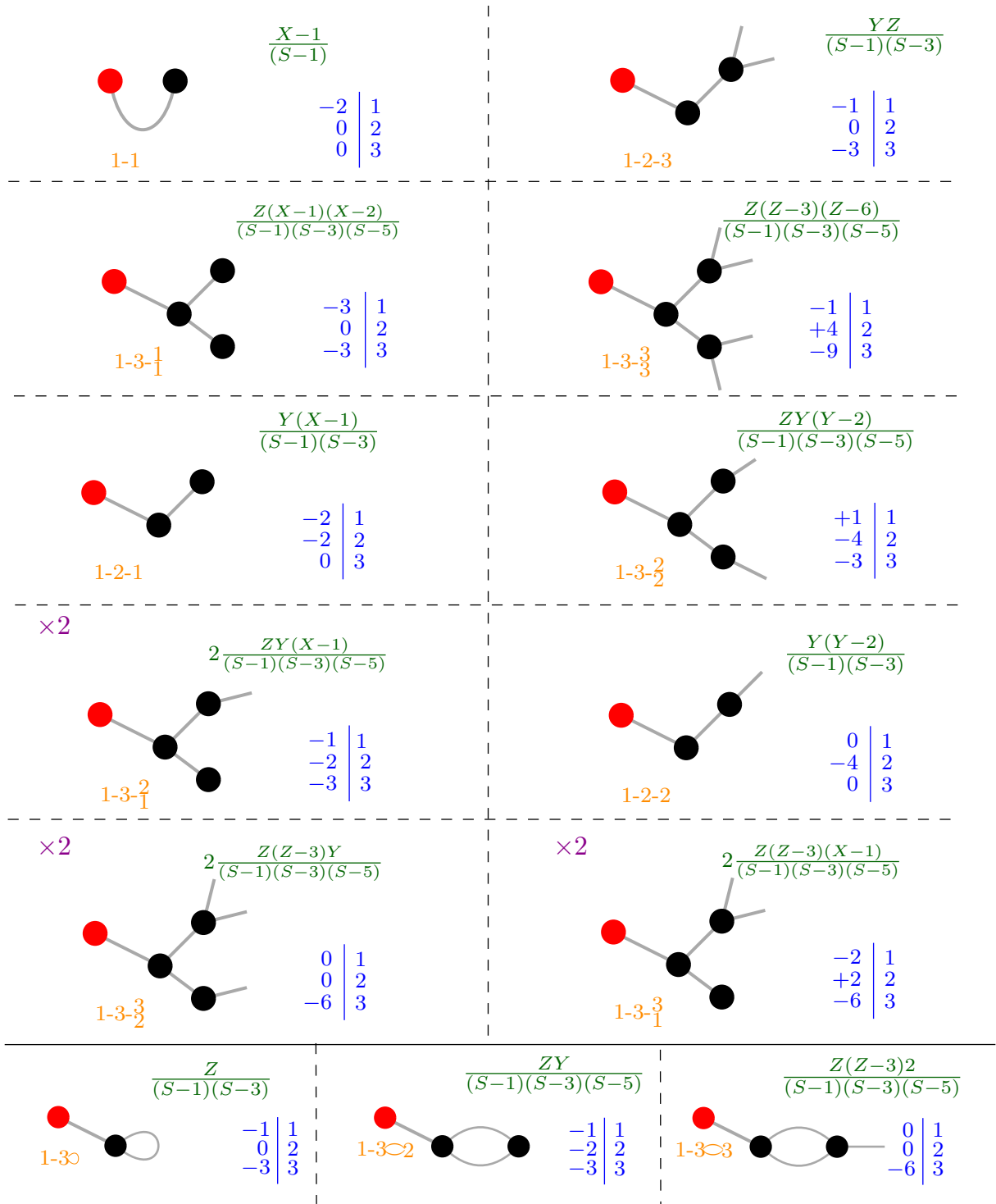


Figure 3: Transitions probabilities of the Markov chain (X^n, Y^n, Z^n) : as long as $X^n > 0$, a vertex ℓ of degree 1 (in red above) is picked and its neighbor v is revealed. The vertices ℓ, v are then removed from the configuration model as well as the connections they created. The probability of each event is indicated in green in the upper right corner. The variation of X, Y, Z are displayed in blue. A symmetry factor is indicated when needed in purple in the upper left corner. Notice in particular that the last three cases on the bottom have probabilities of smaller order $O(1/S)$, so they will not participate to the large scale limit.

In particular, notice that at the stopping time θ^n , the graph made by pairing the remaining unmatched edges is precisely the Karp–Sipser core of $\text{CM}(\mathbf{d}^n)$ and so the second part of Theorem 2 is already proved.

3 Phase transition *via* fluid limit of the Markov chain

In this section, we prove Theorem 1. The main ingredient is a deterministic fluid limit result for the Markov chain (X^n, Y^n, Z^n) .

3.1 Fluid limit for the Markov chain

For a process indexed by discrete time $(\mathfrak{H}_k : k \geq 0)$ we use the notation $\Delta\mathfrak{H}_k = \mathfrak{H}_{k+1} - \mathfrak{H}_k$ for $k \geq 0$. Given the transitions of the Markov chain (X^n, Y^n, Z^n) the following should come as no surprise.

Proposition 2 (Fluid limit). *Suppose that $\mathbf{d}^n = (d_1^n, d_2^n, d_3^n)$ satisfies (1). Then we have the following convergence in probability for the uniform norm:*

$$\left(\frac{X^n_{\lfloor tn \rfloor}}{n}, \frac{Y^n_{\lfloor tn \rfloor}}{n}, \frac{Z^n_{\lfloor tn \rfloor}}{n} \right)_{0 \leq t \leq \theta^n/n} \xrightarrow[n \rightarrow \infty]{(\mathbb{P})} (\mathcal{X}(t), \mathcal{Y}(t), \mathcal{Z}(t))_{0 \leq t \leq t_{\text{ext}}}, \quad (4)$$

where $(\mathcal{X}, \mathcal{Y}, \mathcal{Z})$ is the unique solution⁵ to the differential equation $(\mathcal{X}', \mathcal{Y}', \mathcal{Z}') = \phi(\mathcal{X}, \mathcal{Y}, \mathcal{Z})$ with ϕ defined below (5) with initial conditions (p_1, p_2, p_3) and where t_{ext} is the first hitting time of 0 by the continuous process \mathcal{X} . Moreover, $\theta^n/n \rightarrow t_{\text{ext}}$ in probability as $n \rightarrow \infty$.

Proof. It is a standard application of the differential equation method. Indeed, the increments of the Markov chain (X^n, Y^n, Z^n) are bounded and using the exact transitions (Figure 3), the conditional expected drifts

$$\mathbb{E}[\Delta X_k^n, \Delta Y_k^n, \Delta Z_k^n \mid \mathcal{F}_k]$$

converge for large values of n towards $\phi\left(\frac{X_k^n}{n}, \frac{Y_k^n}{n}, \frac{Z_k^n}{n}\right)$ where the function ϕ is defined by

$$\phi \begin{pmatrix} \mathcal{X} \\ \mathcal{Y} \\ \mathcal{Z} \end{pmatrix} = \begin{pmatrix} -2\mathbf{x} - \mathbf{yz} - 3\mathbf{x}^2\mathbf{z} - 2\mathbf{yx} + \mathbf{zy}^2 - 2\mathbf{zxy} - \mathbf{z}^3 - 4\mathbf{z}^2\mathbf{x} \\ 4\mathbf{z}^3 - 2\mathbf{xy} - 4\mathbf{zy}^2 - 4\mathbf{xyz} - 4\mathbf{y}^2 + 4\mathbf{z}^2\mathbf{x} \\ -3\mathbf{yz} - 3\mathbf{zy}^2 - 12\mathbf{z}^2\mathbf{y} - 3\mathbf{zx}^2 - 6\mathbf{xyz} - 12\mathbf{z}^2\mathbf{x} - 9\mathbf{z}^3 \end{pmatrix}, \quad (5)$$

$$\text{with } \mathcal{S} := \mathcal{X} + \mathcal{Y} + \mathcal{Z} \text{ and where } \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \\ \mathbf{z} \end{pmatrix} := \frac{1}{\mathcal{S}} \begin{pmatrix} \mathcal{X} \\ \mathcal{Y} \\ \mathcal{Z} \end{pmatrix} \text{ is the proportion vector.} \quad (6)$$

For any $\delta > 0$, the convergence of the conditional expected drifts to ϕ is uniform on $\{n^{-1} \cdot S^n \geq \delta\}$ and $(x, y, z) \mapsto \phi(x, y, z)$ is Lipschitz on $\{(x, y, z) \in \mathbb{R}_+^3 : \delta^{-1} \geq x + y + z \geq \delta\}$ as $\nabla\phi(x, y, z)$ is of the form $\frac{P(x, y, z)}{(x+y+z)^4}$, where P is a polynomial. Therefore, by [19, Theorem 1], the equation $(\mathcal{X}', \mathcal{Y}', \mathcal{Z}') = \phi(\mathcal{X}, \mathcal{Y}, \mathcal{Z})$ with initial condition (p_1, p_2, p_3) has a unique solution until the time t_{ext}^δ

⁴The unmatched degree of a vertex at time k is the number of half-edges attached to this vertex which are still unmatched at time k .

⁵More precisely, by *solution*, we mean that $(\mathcal{X}, \mathcal{Y}, \mathcal{Z})$ is a continuous function from $[0, t_{\text{ext}}]$ to \mathbb{R}^3 such that \mathcal{X} first hits 0 at time t_{ext} and $(\mathcal{X}'(t), \mathcal{Y}'(t), \mathcal{Z}'(t)) = \phi(\mathcal{X}(t), \mathcal{Y}(t), \mathcal{Z}(t))$ for all $0 \leq t < t_{\text{ext}}$.

where \mathcal{X} first hits δ , and the convergence (4) holds for $0 \leq t \leq t_{\text{ext}}^\delta$. Moreover, let $t_{\text{ext}} = \lim_{\delta \rightarrow 0} t_{\text{ext}}^\delta$. Since ϕ is bounded by an absolute constant, the solution $(\mathcal{X}, \mathcal{Y}, \mathcal{Z})$ is Lipschitz on $[0, t_{\text{ext}})$, so we can extend it uniquely in a continuous way to $[0, t_{\text{ext}}]$, and by continuity t_{ext} is indeed the first time where \mathcal{X} hits 0. We know that (4) holds on every compact subset of $[0, t_{\text{ext}})$. Moreover, the increments of (X^n, Y^n, Z^n) are bounded by an absolute constant, so the functions $n^{-1} \cdot (X^n, Y^n, Z^n)$ are uniformly Lipschitz and the previous convergence extends to a uniform convergence on $[0, t_{\text{ext}}]$.

We now only need to check that $\frac{\theta^n}{n}$ converges in probability to t_{ext} . We notice that deterministically, if $k < \theta^n$, then $S_{k+1}^n \leq S_k^n - 2$, which implies $\theta^n \leq n$, so up to extraction we may assume that $\frac{\theta^n}{n}$ converges to some random variable \tilde{t}_{ext} . By convergence of the process and the definition of t_{ext} , it is immediate that $\tilde{t}_{\text{ext}} \geq t_{\text{ext}}$. For the other direction, we treat two cases separately:

- if $\mathcal{S}(t_{\text{ext}}) = 0$, then let $\varepsilon > 0$, and let $\delta > 0$ be such that $\mathcal{S}(t_{\text{ext}} - \delta) < \varepsilon$. With probability $1 - o(1)$ as $n \rightarrow +\infty$, we have $S_{\lfloor (t_{\text{ext}} - \delta)n \rfloor}^n < 2\varepsilon n$. Since S^n decreases by at least two at each step, this implies $\theta^n \leq (t_{\text{ext}} - \delta)n + \varepsilon n$, so $\tilde{t}_{\text{ext}} \leq t_{\text{ext}}$.
- if $\mathcal{S}(t_{\text{ext}}) > 0$, we first argue that the first component of $\phi(\mathcal{X}, \mathcal{Y}, \mathcal{Z})$ remains bounded from above by a negative constant along the whole trajectory. Indeed, since \mathcal{S} is bounded from below, we have $\mathcal{Z}' \geq -c\mathcal{Z}$ for some constant c along the trajectory. Hence \mathcal{Z} is bounded from below by a positive constant on $[0, t_{\text{ext}}]$, so \mathbf{y} is bounded away from 1. Since the first component of $\phi(\mathcal{X}, \mathcal{Y}, \mathcal{Z})$ is at most $-\mathbf{y}\mathbf{z} + \mathbf{y}^2\mathbf{z} = -\mathbf{y}\mathbf{z}(1 - \mathbf{y})$, this proves our claim. Therefore, with high probability, the conditional expected drift $\mathbb{E}[\Delta X_k^n | X_k^n]$ is also bounded from above by a negative constant $-c$ along the trajectory. Since the increments are bounded, by the weak law of large numbers this ensures $\tilde{t}_{\text{ext}} \leq t_{\text{ext}}^\varepsilon + \frac{1}{c}\varepsilon$ for all $\varepsilon > 0$, so $\tilde{t}_{\text{ext}} = t_{\text{ext}}$.

□

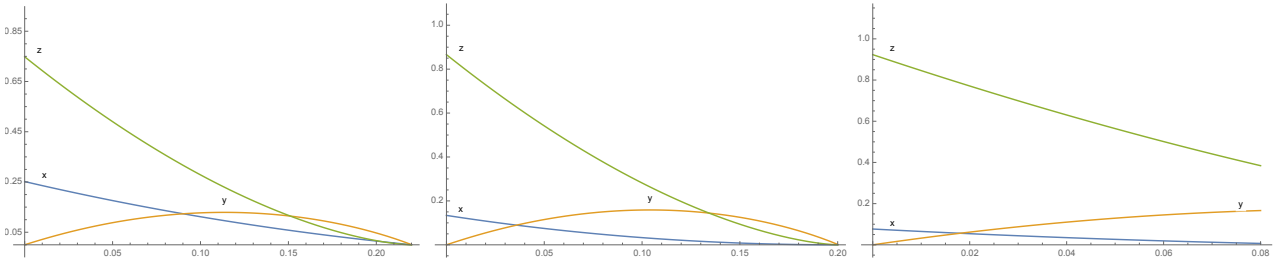


Figure 4: Illustration of the differential system $(\mathcal{X}, \mathcal{Y}, \mathcal{Z})$ in terms of “number of legs” in the subcritical (left), critical (center) and supercritical (right) cases.

3.2 Solving the differential equation

In this section, our goal will be to gather information about the solutions to (5), which will give Theorem 1 and be an important tool in the proof of Theorem 2. As indicated by the system (5), we will see that the solutions are easier to express in terms of proportions. We refer to Figures 4 and 5 for a visualization of the trajectories of these solutions.

Proposition 3. We fix $p_1 > 0$ and $p_2, p_3 \geq 0$ with $p_1 + p_2 + p_3 = 1$. Let $(\mathcal{X}(t), \mathcal{Y}(t), \mathcal{Z}(t))_{0 \leq t \leq t_{\text{ext}}}$ be the solution to (5) with initial condition (p_1, p_2, p_3) . Recall from (2) the definition

$$\Theta = (p_3 - p_1)^2 - 4p_1 \in [-3, 1].$$

- If $\Theta < 0$ (subcritical case), then $\mathcal{X}(t_{\text{ext}}) = \mathcal{Y}(t_{\text{ext}}) = \mathcal{Z}(t_{\text{ext}}) = 0$. Moreover, for $t < t_{\text{ext}}$ sufficiently close to t_{ext} , we have $\mathcal{Z}(t) < \mathcal{X}(t)$.
- If $\Theta > 0$ (supercritical case), then

$$\mathcal{X}(t_{\text{ext}}) = 0, \quad \mathcal{Y}(t_{\text{ext}}) = \frac{4\Theta}{3 + \Theta} (1 - \sqrt{\Theta}) > 0, \quad \text{and} \quad \mathcal{Z}(t_{\text{ext}}) = \frac{4\Theta^{3/2}}{3 + \Theta} > 0. \quad (7)$$

- If $\Theta = 0$ (critical case), then $\mathcal{X}(t_{\text{ext}}) = \mathcal{Y}(t_{\text{ext}}) = \mathcal{Z}(t_{\text{ext}}) = 0$, and more precisely as $\varepsilon \rightarrow 0$:

$$\begin{cases} \mathcal{X}(t_{\text{ext}} - \varepsilon) \sim 3\varepsilon^2, \\ \mathcal{Y}(t_{\text{ext}} - \varepsilon) \sim 4\varepsilon, \\ \mathcal{Z}(t_{\text{ext}} - \varepsilon) \sim 4\sqrt{3}\varepsilon^{3/2}. \end{cases} \quad (8)$$

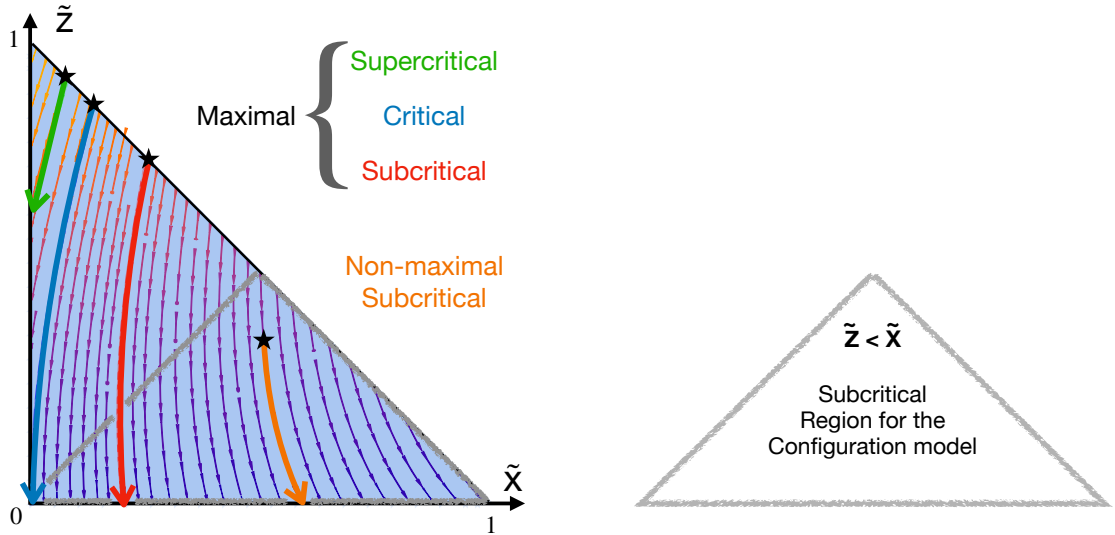


Figure 5: Illustration of the differential system $\tilde{\mathbf{x}}, \tilde{\mathbf{z}}$ with the vector field. The maximal solutions start from $\tilde{\mathbf{x}}(0) + \tilde{\mathbf{z}}(0) = 1$. A maximal supercritical (resp. critical, resp. subcritical) solution is shown in green (resp. blue, resp. red). A non-maximal subcritical solution is displayed in orange. Note that any subcritical solution terminates in the gray region which is subcritical for the configuration model itself.

Proof. We will first obtain an explicit (up to time-change) solution to (5). We recall that $\mathcal{S} = \mathcal{X} + \mathcal{Y} + \mathcal{Z}$ is the fluid limit of the sum process and that $\mathbf{x}, \mathbf{y}, \mathbf{z}$ are the proportions whose sum is constant and equal to 1.

Using $\mathbf{y} = 1 - \mathbf{x} - \mathbf{z}$, the system (5) translates into the following system on \mathbf{x}, \mathbf{z} and \mathcal{S} :

$$\begin{cases} \mathbf{x}' = \frac{1}{\mathcal{S}}(\mathbf{x} - \mathbf{z})\mathbf{z}, \\ \mathbf{z}' = \frac{1}{\mathcal{S}}(-2 + \mathbf{x} - \mathbf{z})\mathbf{z}, \\ \mathcal{S}' = 2(-2 + \mathbf{x} - \mathbf{z}), \end{cases} \quad (9)$$

where again $\mathcal{S}(0) = 1$ and $\mathbf{x}(0), \mathbf{z}(0) \geq 0$ satisfy $\mathbf{x}(0) + \mathbf{z}(0) \leq 1$.

In order to get rid of \mathcal{S} in this system, we perform a time change: for $t \in [0, t_{\text{ext}}]$, we write

$$\gamma(t) = \int_0^t \frac{ds}{\mathcal{S}(s)} \in [0, +\infty].$$

We also define the functions $\tilde{\mathbf{x}}, \tilde{\mathbf{y}}, \tilde{\mathbf{z}}$ on $[0, u_{\text{ext}}]$, with $u_{\text{ext}} = \int_0^{t_{\text{ext}}} \frac{ds}{\mathcal{S}(s)}$, by $\tilde{\mathbf{x}}(u) = \mathbf{x}(\gamma^{-1}(u))$ and $\tilde{\mathbf{z}}(u) = \mathbf{z}(\gamma^{-1}(u))$. We obtain the system

$$\begin{cases} \tilde{\mathbf{x}}' &= (\tilde{\mathbf{x}} - \tilde{\mathbf{z}})\tilde{\mathbf{z}}, \\ \tilde{\mathbf{z}}' &= (-2 + \tilde{\mathbf{x}} - \tilde{\mathbf{z}})\tilde{\mathbf{z}}. \end{cases}$$

We find solutions to this system as follows: by subtracting the second line to the first one, we have $\tilde{\mathbf{x}}' - \tilde{\mathbf{z}}' = 2\tilde{\mathbf{z}}$ and the second line implies that $\tilde{\mathbf{x}} - \tilde{\mathbf{z}} = \left(\frac{\tilde{\mathbf{z}}'}{\tilde{\mathbf{z}}} + 2\right)$. Deriving the second identity and comparing, we deduce the following second-order non-linear one-dimensional differential equation:

$$2(\tilde{\mathbf{z}})^3 = \tilde{\mathbf{z}}''\tilde{\mathbf{z}} - (\tilde{\mathbf{z}}')^2.$$

A complete family of solutions is given by

$$\begin{cases} \tilde{\mathbf{z}}(u) &= \frac{b^2}{\sinh(b(u + u_0))^2}, \\ \tilde{\mathbf{x}}(u) &= \left(\frac{b}{\tanh(b(u + u_0))} - 1\right)^2 + 1 - b^2, \\ \tilde{\mathbf{y}}(u) &= \frac{-2b^2}{\tanh^2(b(u + u_0))} + \frac{2b}{\tanh(b(u + u_0))} + 2b^2 - 1, \end{cases} \quad (10)$$

where $b, u_0 \in \mathbb{R}$. We notice that along these solutions, the quantity $(\tilde{\mathbf{z}} - \tilde{\mathbf{x}})^2 - 4\tilde{\mathbf{x}}$ is constant, and is equal to $4(b^2 - 1)$, this quantity is equal to the Θ defined by (2):

$$(\tilde{\mathbf{z}} - \tilde{\mathbf{x}})^2 - 4\tilde{\mathbf{x}} \equiv 4(b^2 - 1) = (p_3 - p_1)^2 - 4p_1 = \Theta. \quad (11)$$

We also notice that $\tilde{\mathbf{y}}$ is always increasing and that $\tilde{\mathbf{y}} < 0$ for u small enough, which has no meaning in our context. Therefore, every solution is contained in a *maximal* solution, i.e. a solution where the initial condition (p_1, p_2, p_3) satisfies $p_2 = 0$. Since we are only interested in the behavior near extinction and since the right-hand side of the formulas (7) depends only on b , we may restrict ourselves to maximal solutions, i.e. assume $p_2 = 0$. Other solutions can be deduced from this by a time shift in $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}, \tilde{\mathbf{z}})$, which translates into a time shift in $(\mathcal{X}, \mathcal{Y}, \mathcal{Z})$. From $\tilde{\mathbf{y}}(0) = 0$, we get

$$u_0 = \frac{1}{2b} \log \left(1 + 2b + \frac{2b\sqrt{(4b^2 - 1)}}{2b - 1} \right) > 0,$$

so $p_1 = 1 - \frac{1}{2}\sqrt{4b^2 - 1}$ and $p_3 = \frac{1}{2}\sqrt{4b^2 - 1}$.

We now come back to the true solutions $(\mathcal{X}, \mathcal{Y}, \mathcal{Z})$ in each of the three cases of Proposition 3. For this, we need to study the time change $\gamma : [0, t_{\text{ext}}] \rightarrow [0, u_{\text{ext}}]$. By definition of γ and the third line of (9), for all $t \in [0, t_{\text{ext}})$, we have

$$\begin{cases} \frac{1}{\mathcal{S}(t)} &= \gamma'(t), \\ \mathcal{S}'(t) &= 2(-2 + \tilde{\mathbf{x}}(\gamma(t)) - \tilde{\mathbf{z}}(\gamma(t))). \end{cases}$$

Multiplying those lines and integrating both sides using the exact expressions of $\tilde{\mathbf{x}}$ and $\tilde{\mathbf{z}}$, we find $\frac{d}{dt} \log \mathcal{S}(t) = -4 \frac{d}{dt} \log(\sinh(b \cdot (\gamma(t) + u_0)))$ so the following quantity is constant:

$$\mathcal{S}(t) \sinh^4(b \cdot (\gamma(t) + u_0)) = \mathcal{S}(t) \left(\frac{b^2}{\mathbf{z}(t)} \right)^2 = \frac{b^4}{\tilde{\mathbf{z}}(0)^2} = \frac{4b^4}{4b^2 - 1}. \quad (12)$$

Note that this last equation, combined with the expression of $\mathcal{S}'(t)$, provides a differential equation satisfied by \mathcal{S} , from which we could express \mathcal{S} as the inverse bijection of an explicit function. However, this will not be needed in the proof. Given those findings, the rest of the proof is made of easy calculations. Let us proceed. We refer to Figures 4 and 5 for visualization of the system in terms of proportions or in “number of legs”.

Subcritical regime. For $\Theta < 0$, we have $\frac{1}{2} < b < 1$. In this case, we observe that $\tilde{\mathbf{x}}(u) \geq 1 - b^2$ is bounded away from 0, so the same is true for $\mathbf{x}(t)$ on $[0, t_{\text{ext}}]$. It follows that $\mathcal{S}(t_{\text{ext}}) = \frac{\mathcal{X}(t_{\text{ext}})}{\mathbf{x}(t_{\text{ext}})} = 0$. Therefore, by (12), we have $\mathbf{z}(t_{\text{ext}}) = \sqrt{\left(\frac{4b^2-1}{4}\right) \mathcal{S}(t_{\text{ext}})} = 0$. In particular, for t sufficiently close to t_{ext} , we have $\mathbf{z}(t) < \mathbf{x}(t)$, so $\mathcal{Z}(t) < \mathcal{X}(t)$. Note that this also implies $u_{\text{ext}} = +\infty$.

Supercritical regime. For $\Theta > 0$, we have $1 < b < \frac{\sqrt{5}}{2}$. In this case, the function $\tilde{\mathbf{x}}$ first hits 0 at time

$$\hat{u}_{\text{ext}} = -u_0 + \frac{1}{b} \operatorname{Arccoth} \frac{1 + \sqrt{b^2 - 1}}{b}.$$

This implies that $u_{\text{ext}} \leq \hat{u}_{\text{ext}}$. We claim that we have equality. Indeed, if this is not the case, we have $\mathbf{x}(t_{\text{ext}}) = \tilde{\mathbf{x}}(u_{\text{ext}}) > 0$, so $\mathcal{S}(t_{\text{ext}}) = 0$, so (12) implies $\tilde{\mathbf{z}}(u_{\text{ext}}) = 0$ with $u_{\text{ext}} < +\infty$, which is not possible given the explicit expression of $\tilde{\mathbf{z}}$. Therefore, we have $\tilde{\mathbf{x}}(u_{\text{ext}}) = 0$. Using (11) we can compute

$$\mathbf{z}(t_{\text{ext}}) = \tilde{\mathbf{z}}(u_{\text{ext}}) = 2\sqrt{b^2 - 1} \quad \text{and} \quad \mathbf{y}(t_{\text{ext}}) = 1 - 2\sqrt{b^2 - 1}$$

and finally, using (12):

$$\mathcal{S}(t_{\text{ext}}) = \frac{4}{4b^2 - 1} \tilde{\mathbf{z}}(u_{\text{ext}}) = \frac{16(b^2 - 1)}{4b^2 - 1},$$

which, once translated in terms of Θ , gives (7).

Critical regime. For $\Theta = 0$, the maximal solution starts from $p_3 = \frac{\sqrt{3}}{2}$ and $p_1 = 1 - \frac{\sqrt{3}}{2}$, and we have $b = 1$. In particular, using $u_0 \geq 0$, we have $\tilde{\mathbf{x}}(u) > 0$ for all $u \geq 0$. By the same argument as in the supercritical regime, this implies $u_{\text{ext}} = +\infty$. Therefore, by the exact expression of $\tilde{\mathbf{x}}, \tilde{\mathbf{y}}, \tilde{\mathbf{z}}$, as $t \rightarrow t_{\text{ext}}$, we have $\mathbf{x}(t), \mathbf{z}(t) \rightarrow 0$ and $\mathbf{y}(t) \rightarrow 1$. Therefore, by (12) at $t = t_{\text{ext}}$, we have $\mathcal{S}(t_{\text{ext}}) = 0$, so $\mathcal{Y}(t_{\text{ext}}) = \mathcal{Z}(t_{\text{ext}}) = 0$.

Hence, letting $t \rightarrow t_{\text{ext}}$ in the third equation of (9), we have $\mathcal{S}'(t) \rightarrow -4$ as $t \rightarrow t_{\text{ext}}$, so $\mathcal{S}(t_{\text{ext}} - \varepsilon) \sim 4\varepsilon$ as $\varepsilon \rightarrow 0$. Injecting this in (12), we find $\mathbf{z}(t_{\text{ext}} - \varepsilon) \sim \sqrt{3}\varepsilon$, so $\mathcal{Z}(t_{\text{ext}} - \varepsilon) \sim 4\sqrt{3}\varepsilon^{3/2}$. Finally, we know from (11) that $(\mathbf{z} - \mathbf{x})^2 - 4\mathbf{x}$ is constant equal to 0, so $\mathbf{x}(t_{\text{ext}} - \varepsilon) \sim \frac{1}{4}\mathbf{z}(t_{\text{ext}} - \varepsilon)^2 \sim \frac{3}{4}\varepsilon$, which gives the asymptotics for \mathcal{X} . \square

3.3 Phase transition: proof of Theorem 1

Subcritical regime. We assume that (p_1, p_2, p_3) is subcritical, and consider the associated solution $(\mathcal{X}, \mathcal{Y}, \mathcal{Z})$ to the differential equation. By Proposition 3, let $t_1 < t_{\text{ext}}$ be such that $\mathcal{Z}(t_1) < \mathcal{X}(t_1)$.

By Proposition 2, we have

$$\frac{1}{n} \left(X_{[t_1 n]}^n, Y_{[t_1 n]}^n, Z_{[t_1 n]}^n \right) \xrightarrow[n \rightarrow +\infty]{(\mathbb{P})} (\mathcal{X}(t_1), \mathcal{Y}(t_1), \mathcal{Z}(t_1)).$$

Moreover, by Proposition 1, conditionally on $\mathcal{F}_{[t_1 n]}$, the remaining graph after $[t_1 n]$ steps of the Karp–Sipser algorithm is a configuration model with respectively $X_{[t_1 n]}^n$, $Y_{[t_1 n]}^n$ and $Z_{[t_1 n]}^n$ half-edges belonging to vertices of degree 1, 2 and 3. Since $n^{-1} Z_{[t_1 n]}^n \approx \mathcal{Z}(t_1) < \mathcal{X}(t_1) \approx n^{-1} X_{[t_1 n]}^n$ this is a *subcritical* configuration model (do not confuse with subcriticality in terms of the Karp–Sipser core). In particular, by [15, Theorem 1.b] there is a constant $c = c(p_1, p_2, p_3)$ such that with high probability the remaining subgraph after $[t_1 n]$ steps has fewer than $c \log(n)$ cycles and all of its connected components have size at most $c \log(n)$. On the other hand, by construction, the Karp–Sipser core is included in the union of all the cycles of $G_{[t_1 n]}^n$, so it has size $O_{\mathbb{P}}(\log^2 n)$.

Remark (True size of the subcritical KS-core). The above bound $O_{\mathbb{P}}(\log^2 n)$ for the size of the subcritical Karp–Sipser core is very crude towards the end of the proof. We expect the actual order of magnitude of the KS-core to be $O_{\mathbb{P}}(1)$ as in the Erdős–Rényi case [1].

Critical and supercritical regime. In this case, combining Proposition 2 and our explicit computations of the solutions, we obtain that $(X^n/S^n, Y^n/S^n, Z^n/S^n, n^{-1} \cdot S^n)(\theta^n)$ converges to

$$(\mathbf{x}(t_{\text{ext}}), \mathbf{y}(t_{\text{ext}}), \mathbf{z}(t_{\text{ext}}), \mathcal{S}(t_{\text{ext}})) = \left(0, 1 - 2\sqrt{b^2 - 1}, 2\sqrt{b^2 - 1}, \frac{16(b^2 - 1)}{4b^2 - 1} \right).$$

In particular the number of half-edges of the Karp–Sipser core is equal to $S_{\theta^n}^n = Y_{\theta^n}^n + Z_{\theta^n}^n$, so it is asymptotically $o_{\mathbb{P}}(n)$ if $b = 1$ (critical case). If $b > 1$, it is linear in n , which concludes the proof of Theorem 1 after a quick computation.

4 Analysis of the critical case

In this section, we shall prove our main result Theorem 2. In the rest of the paper, we shall thus suppose that the initial conditions (3) are in force. Let us first explain the heuristics to help the reader follow the proof. We refer to Figure 6 for an illustration.

We have seen above that in the critical regime, the asymptotic size of the Karp–Sipser core is $o_{\mathbb{P}}(n)$ and that almost all vertices have degree 2 (i.e. with density 1 since $\mathbf{y}(t_{\text{ext}}) = 1$). Recall that the process stops at time

$$\theta^n = \inf\{k \geq 0 : X_k^n = 0\},$$

which by Proposition 2 is $\approx t_{\text{ext}} \cdot n$. To analyse this stopping time and understand the size of the KS-core, we need to be more precise in the analysis of the fluctuations of the process (X^n, Y^n, Z^n) around its fluid limit $n \cdot (\mathcal{X}, \mathcal{Y}, \mathcal{Z})$. To this end, we define the fluctuations processes $(A_k^n, B_k^n, C_k^n)_{0 \leq k \leq \theta^n}$ by

$$\begin{cases} X_k^n &= n\mathcal{X}\left(\frac{k}{n}\right) + A_k^n \\ Y_k^n &= n\mathcal{Y}\left(\frac{k}{n}\right) + B_k^n \\ Z_k^n &= n\mathcal{Z}\left(\frac{k}{n}\right) + C_k^n \end{cases}$$

To simplify notation, the n in the exponent will be implicit for the rest of the paper when there is no ambiguity, even if we will often look at the asymptotic as n goes to infinity.

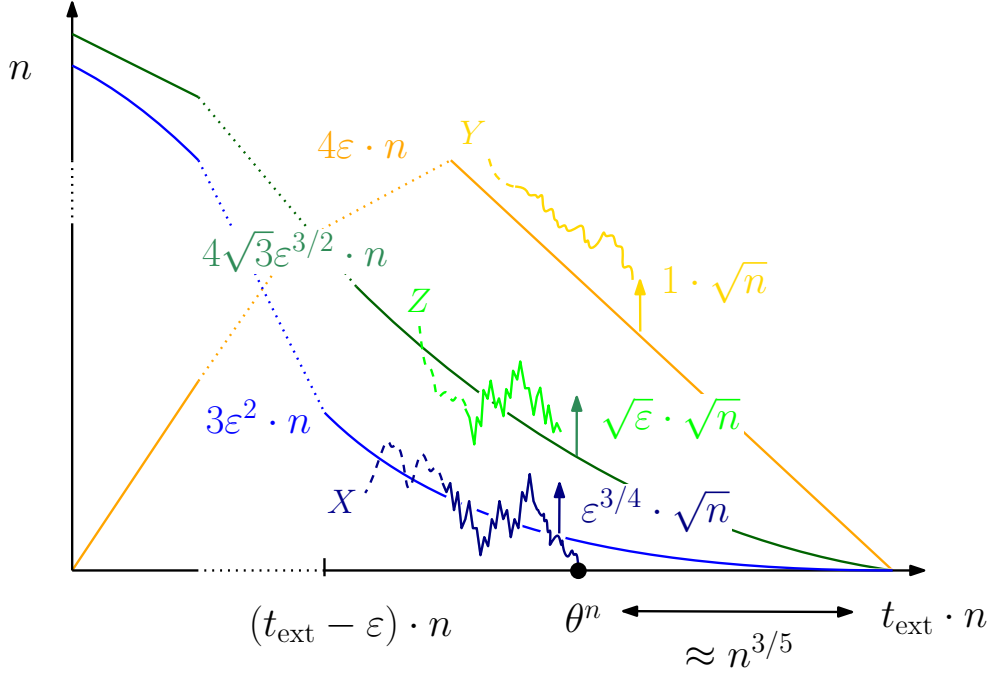


Figure 6: Heuristics for the proof of Theorem 2. The variations of the processes (X, Y, Z) around its deterministic fluid limit when $k = (t_{\text{ext}} - \varepsilon_k)n$ are displayed above. In particular, in the case of X , the number of degree 1 vertices, those variations may cause X to touch 0 when $\varepsilon_k \approx n^{-2/5}$ so that there are $\varepsilon_k n \approx n^{3/5}$ vertices of degree 2 and $\varepsilon_k^{3/2} n \approx n^{2/5}$ vertices of degree 3 remaining in the graph.

When we are sufficiently far from the end of the process, i.e. when $k \approx tn$ for $0 \leq t < t_{\text{ext}}$ we know from Proposition 2 that (X, Y, Z) is well approximated by $n \cdot (\mathcal{X}, \mathcal{Y}, \mathcal{Z})$ and classical results (see Lemma 2) will show that the fluctuations A, B and C renormalized by a factor $1/\sqrt{n}$ converge to Gaussian variables whose variances depend on t . To analyse the algorithm towards the end we will use the notation, for $0 \leq k \leq (t_{\text{ext}}n) \wedge \theta$,

$$\boxed{\varepsilon_k := t_{\text{ext}} - \frac{k}{n} \geq 0 \quad \text{so that} \quad k = (t_{\text{ext}} - \varepsilon_k)n.} \quad (13)$$

Notice the bold font for ε to avoid confusion. Recall from Equation (8) that $\mathcal{X}(k/n), \mathcal{Y}(k/n)$ and $\mathcal{Z}(k/n)$ are of order respectively $\varepsilon_k^2, \varepsilon_k$ and $\varepsilon_k^{3/2}$. We will see below that the order of magnitude of $n^{-1/2} \cdot A_k, n^{-1/2} \cdot B_k$ and $n^{-1/2} \cdot C_k$ are respectively $\varepsilon_k^{3/4}, 1$ and $\varepsilon_k^{1/2}$. In particular, the fluctuations A of X become of the same order of magnitude as its deterministic approximation $n\mathcal{X}$ when

$$n\varepsilon_k^2 \approx n\mathcal{X}(t_{\text{ext}} - \varepsilon_k) \approx A_k \approx \sqrt{n} \cdot \varepsilon_k^{3/4} \quad \text{i.e. when} \quad \varepsilon_k \approx n^{-2/5} \iff nt_{\text{ext}} - k \approx n^{3/5},$$

and this explains heuristically why $\theta_n = t_{\text{ext}}n + O(n^{3/5})$ and why the size of the Karp-Sipser core is given essentially by $Y_{\theta_n} \approx n^{3/5}$. The rest of this section makes those heuristic rigorous and proves our main result Theorem 2.

We first provide estimations of the conditional expected drifts and variances of the increments of the fluctuation processes (A, B, C) in Propositions 4 and 5. These propositions support the above heuristics and lead us to introduce the renormalized fluctuations processes

$$\tilde{A}_k = \frac{A_k}{\varepsilon_k^{3/4} \sqrt{n}}, \quad \tilde{B}_k = \frac{B_k}{\sqrt{n}}, \quad \text{and} \quad \tilde{C}_k = \frac{C_k}{\varepsilon_k^{1/2} \sqrt{n}},$$

which, at least heuristically, should be tight in k . After that, our proof consists in two main steps. First we will show that with high probability as $n \rightarrow \infty$, we can bound –with some log’s– the process $(\tilde{A}, \tilde{B}, \tilde{C})$ up to time $O(n^{3/5})$ before $t_{\text{ext}}n$, see Proposition 6. To do so we will extensively use the fact that for \tilde{C} and \tilde{A} , the conditional expected drifts tend “to pull them back to 0” so that the processes remain small over all scales. Finally, in a second step, we will show that when $k = nt_{\text{ext}} - tn^{3/5}$ for $x \in \mathbb{R}$ the fluctuation process \tilde{A} is well approximated by a stochastic differential equation, see Proposition 8. The fluctuations \tilde{B} and \tilde{C} are, at this scale, still negligible in front of their differential method approximation.

4.1 Drift and variance estimates

In this section we compute the conditional expected drift and variance of the fluctuations processes A, B, C . Recall the very important notation ε_k introduced in (13). As explained above, it will turn out that $\theta \equiv \theta^n$ is located around $t_{\text{ext}}n - O(n^{3/5})$ and in the forthcoming Propositions 4 and 5 we shall allow a little room and only look at times $k < \theta$ such that $\varepsilon_k \geq n^{-2/5-1/100}$ (and indeed the fraction $1/100$ is somehow arbitrary). We thus put

$$\tilde{\theta} = \theta \wedge \left(t_{\text{ext}}n - n^{3/5-1/100} \right). \quad (14)$$

Recall from above the notation

$$\tilde{A}_k := \frac{X_k - n\mathcal{X}\left(\frac{k}{n}\right)}{\varepsilon_k^{3/4}\sqrt{n}} \quad \tilde{B}_k := \frac{Y_k - n\mathcal{Y}\left(\frac{k}{n}\right)}{\sqrt{n}} \quad \text{and} \quad \tilde{C}_k := \frac{Z_k - n\mathcal{Z}\left(\frac{k}{n}\right)}{\varepsilon_k^{1/2}\sqrt{n}}.$$

Recall also that \mathcal{F}_k is the σ -algebra generated by $(X_i, Y_i, Z_i)_{0 \leq i \leq k}$. We have chosen the normalization so that the processes \tilde{A}_k, \tilde{B}_k and \tilde{C}_k are of order 1 and fluctuate at the time-scale $\varepsilon_k n$, which is why the conditional expected drift and variances are all of order $\frac{1}{\varepsilon_k n}$.

Proposition 4 (Drift estimates). *There exists a constant $K > 0$ such that for all $\delta > 0$, there is $\eta \equiv \eta(\delta) > 0$ such that the following holds for n large enough. For any $(t_{\text{ext}} - \eta)n \leq k < \tilde{\theta}$, if we have $|\tilde{A}_k|, |\tilde{B}_k|, |\tilde{C}_k| < 1000 \log n$ then:*

$$\left| \mathbb{E} \left[\Delta \tilde{A}_k | \mathcal{F}_k \right] + \frac{1}{\varepsilon_k n} \frac{1}{4} \tilde{A}_k \right| \leq \frac{\delta}{\varepsilon_k n} |\tilde{A}_k| + \frac{K \varepsilon_k^{1/4}}{\varepsilon_k n} \max \left(|\tilde{B}_k|, |\tilde{C}_k| \right) + \frac{K}{\varepsilon_k n} n^{-1/30}, \quad (15)$$

$$\left| \mathbb{E} \left[\Delta \tilde{B}_k | \mathcal{F}_k \right] \right| \leq \frac{K}{\varepsilon_k n} \sqrt{\varepsilon_k} \max \left(|\tilde{A}_k|, |\tilde{B}_k|, |\tilde{C}_k| \right) + \frac{K}{\varepsilon_k n} n^{-1/30}, \quad (16)$$

$$\left| \mathbb{E} \left[\Delta \tilde{C}_k | \mathcal{F}_k \right] - \frac{1}{\varepsilon_k n} \left(\frac{3\sqrt{3}}{2} \tilde{B}_k - \tilde{C}_k \right) \right| \leq \frac{\delta}{\varepsilon_k n} \max \left(|\tilde{B}_k|, |\tilde{C}_k| \right) + \frac{K}{\varepsilon_k n} \varepsilon_k^{3/4} |\tilde{A}_k| + \frac{K}{\varepsilon_k n} n^{-1/30}. \quad (17)$$

Proposition 5 (Variance estimates). *There exists a constant K such that for all $\delta > 0$, there is $\eta \equiv \eta(\delta) > 0$ such that the following holds for n large enough. For any $(t_{\text{ext}} - \eta)n \leq k < \tilde{\theta}$, if we have $|\tilde{A}_k|, |\tilde{B}_k|, |\tilde{C}_k| < 1000 \log n$ then:*

$$\left| \text{Var} \left(\Delta \tilde{A}_k | \mathcal{F}_k \right) - \frac{2\sqrt{3}}{\varepsilon_k n} \right| \leq \frac{\delta}{\varepsilon_k n} + \frac{K}{\varepsilon_k n} n^{-1/30} + \frac{K \varepsilon_k^{1/2}}{\varepsilon_k n} \tilde{A}_k^2 + \frac{K}{n} \max \left(\tilde{B}_k^2, \tilde{C}_k^2 \right), \quad (18)$$

$$\text{Var} \left(\Delta \tilde{A}_k | \mathcal{F}_k \right) \leq \frac{2\sqrt{3} + \delta}{\varepsilon_k n} + \frac{K}{\varepsilon_k n} n^{-1/30}, \quad (19)$$

$$\text{Var} \left(\Delta \tilde{B}_k | \mathcal{F}_k \right) \leq \frac{K \varepsilon_k}{\varepsilon_k n} \quad (20)$$

$$\text{Var} \left(\Delta \tilde{C}_k | \mathcal{F}_k \right) \leq \frac{K \varepsilon_k^{1/2}}{\varepsilon_k n}. \quad (21)$$

The proofs of the above two propositions follow by examining precisely the probability transitions of the Markov chain (X, Y, Z) given by Figure 3 and basic (though important) analysis of the behavior of the function ϕ (defined by (5)) and its gradient $\nabla\phi$ near t_{ext} . Let us start with a deterministic lemma based on (8) controlling X, Y, Z from the processes $\tilde{A}, \tilde{B}, \tilde{C}$:

Lemma 1. *There are absolute constants $C, c > 0$ such that if $|\tilde{A}_k|, |\tilde{B}_k|, |\tilde{C}_k| < 1000 \log n$ and $X_k > 0$ and $\varepsilon_k \in [n^{-2/5-1/100}, \eta]$, for n large enough we have*

$$X_k \leq C\varepsilon_k^2 n \times n^{1/100}, \quad Y_k \leq C\varepsilon_k n, \quad Z_k \leq C\varepsilon_k^{3/2} n$$

and

$$S_k \geq Y_k \geq c\varepsilon_k n.$$

Proof. Recall the asymptotics (8). We simply write

$$X_k \leq n\mathcal{X}\left(\frac{k}{n}\right) + \varepsilon_k^{3/4} \sqrt{n}\tilde{A}_k \stackrel{(8)}{\leq} C'\varepsilon_k^2 n + 1000\varepsilon_k^{3/4} \sqrt{n} \log n.$$

The assumption $k \leq t_{\text{ext}}n - n^{3/5-1/100}$, i.e. $\varepsilon_k \geq n^{-2/5-1/100}$, implies that the second term is $O(\varepsilon_k^2 n \times n^{1/100})$. The other two upper bounds can be proved in the same way. Finally, we have

$$S_k \geq Y_k = n\mathcal{Y}(k/n) + \sqrt{n}\tilde{B}_k \geq c'\varepsilon_k n - 1000\sqrt{n} \log n,$$

which is enough to prove the lower bound on S_k since $\varepsilon_k n \geq n^{3/5-1/100}$ is much larger than $\sqrt{n} \log n$ if n is large enough. \square

Proof of Proposition 4. Recall the definition of ϕ in (5) given in terms of proportions, so that using the notation $s = x + y + z$ we have

$$\begin{aligned} \phi_X(x, y, z) &= -2\frac{x}{s} - \frac{yz}{s^2} - 3\frac{x^2z}{s^3} - 2\frac{xy}{s^2} + \frac{y^2z}{s^3} - 2\frac{xyz}{s^3} - \frac{z^3}{s^3} - 4\frac{xz^2}{s^3}, \\ \phi_Y(x, y, z) &= 2\left(2\frac{z^3}{s^3} - \frac{xy}{s^2} - 2\frac{y^2z}{s^3} - 2\frac{xyz}{s^3} - 2\frac{y^2}{s^2} + 2\frac{xz^2}{s^3}\right), \\ \phi_Z(x, y, z) &= 3\left(-\frac{yz}{s^2} - \frac{y^2z}{s^3} - 4\frac{yz^2}{s^3} - \frac{x^2z}{s^3} - 2\frac{xyz}{s^3} - 4\frac{xz^2}{s^3} - 3\frac{z^3}{s^3}\right), \end{aligned}$$

and the fluid limit equation is $\mathcal{X}' = \phi_X(\mathcal{X}, \mathcal{Y}, \mathcal{Z})$, and similarly for the two other coordinates.

WE START WITH THE ESTIMATE (15) ON \tilde{A} . We first decompose the conditional expected drift as follows:

$$\begin{aligned} \mathbb{E}\left[\Delta\tilde{A}_k|\mathcal{F}_k\right] &= \frac{1}{\varepsilon_{k+1}^{3/4}\sqrt{n}}\mathbb{E}\left[\Delta A_k|\mathcal{F}_k\right] + \left(\frac{\varepsilon_k^{3/4}}{\varepsilon_{k+1}^{3/4}} - 1\right)\tilde{A}_k \\ &= \frac{1}{\varepsilon_{k+1}^{3/4}\sqrt{n}}\left(\mathbb{E}\left[\Delta X_k|\mathcal{F}_k\right] - n\left(\mathcal{X}\left(\frac{k+1}{n}\right) - \mathcal{X}\left(\frac{k}{n}\right)\right)\right) + \left(\frac{\varepsilon_k^{3/4}}{\varepsilon_{k+1}^{3/4}} - 1\right)\tilde{A}_k \end{aligned}$$

Therefore, by decomposing $1/4 = 1 - 3/4$, we can decompose the left-hand side of (15) as follows:

$$\begin{aligned} & \left| \mathbb{E} \left[\Delta \tilde{A}_k | \mathcal{F}_k \right] + \frac{1}{\varepsilon_k n} \frac{1}{4} \tilde{A}_k \right| \\ & \leq \left| \left(\frac{\varepsilon_k^{3/4}}{\varepsilon_{k+1}^{3/4}} - 1 \right) \tilde{A}_k - \frac{3}{4} \frac{1}{\varepsilon_k n} \tilde{A}_k \right| \end{aligned} \quad (22)$$

$$+ \frac{1}{\varepsilon_{k+1}^{3/4} \sqrt{n}} \left| \mathbb{E} [\Delta X_k | \mathcal{F}_k] - \phi_X \left(\frac{X_k}{n}, \frac{Y_k}{n}, \frac{Z_k}{n} \right) \right| \quad (23)$$

$$+ \frac{1}{\varepsilon_{k+1}^{3/4} \sqrt{n}} \left| \phi_X \left(\frac{X_k}{n}, \frac{Y_k}{n}, \frac{Z_k}{n} \right) - \phi_X \left((\mathcal{X}, \mathcal{Y}, \mathcal{Z}) \left(\frac{k}{n} \right) \right) - \left(\frac{A_k}{n}, \frac{B_k}{n}, \frac{C_k}{n} \right) \cdot \nabla \phi_X \left((\mathcal{X}, \mathcal{Y}, \mathcal{Z}) \left(\frac{k}{n} \right) \right) \right| \quad (24)$$

$$+ \frac{1}{\varepsilon_{k+1}^{3/4} \sqrt{n}} \left| \left(\frac{A_k}{n}, \frac{B_k}{n}, \frac{C_k}{n} \right) \cdot \nabla \phi_X \left((\mathcal{X}, \mathcal{Y}, \mathcal{Z}) \left(\frac{k}{n} \right) \right) + \frac{1}{\varepsilon_k} \frac{A_k}{n} \right| \quad (25)$$

$$+ \frac{1}{\varepsilon_{k+1}^{3/4} \sqrt{n}} \left| \phi_X \left((\mathcal{X}, \mathcal{Y}, \mathcal{Z}) \left(\frac{k}{n} \right) \right) - n \left(\mathcal{X} \left(\frac{k+1}{n} \right) - \mathcal{X} \left(\frac{k}{n} \right) \right) \right|. \quad (26)$$

We will bound each of these five error terms one by one. More precisely, we will prove that the terms (22), (23), (24) and (26) are all $O\left(\frac{n^{-1/30}}{\varepsilon_k n}\right)$, whereas the other terms in (15) come from (25). We start with (22), which is easy. We simply write $\varepsilon_k = t_{\text{ext}} - \frac{k}{n}$ and $\varepsilon_{k+1} = t_{\text{ext}} - \frac{k+1}{n}$. This implies $\frac{\varepsilon_{k+1}}{\varepsilon_k} = 1 - \frac{1}{\varepsilon_k n}$, so

$$\frac{\varepsilon_k^{3/4}}{\varepsilon_{k+1}^{3/4}} - 1 = \frac{3}{4} \frac{1}{\varepsilon_k n} + O\left(\frac{1}{(\varepsilon_k n)^2}\right),$$

where the constant is absolute. Finally, using $\varepsilon_k n \geq n^{3/5-1/100}$ we have

$$\frac{|\tilde{A}_k|}{(\varepsilon_k n)^2} \leq \frac{100 n^{-3/5+1/100} \log n}{\varepsilon_k n},$$

so we can bound (22) by $\frac{K}{\varepsilon_k n} n^{-1/30}$.

We now move on to (23). The drift $\mathbb{E}[\Delta X_k | \mathcal{F}_k]$ can be expressed as the sum over all the cases of Figure 3 of the probability of each case multiplied by the variation of X in this case. For example, the probability for the first case is $\frac{X_k-1}{S_k-1}$. Approximating $\mathbb{E}[\Delta X_k | \mathcal{F}_k]$ by $\phi_X\left(\frac{X_k}{n}, \frac{Y_k}{n}, \frac{Z_k}{n}\right)$ is then equivalent to approximating $\frac{X_k-1}{S_k-1}$ by $\frac{X_k/n}{S_k/n}$, and similarly for all the other terms. But we have

$$\frac{X_k-1}{S_k-1} = \frac{X_k/n}{S_k/n} \times \frac{1 - \frac{1}{X_k}}{1 - \frac{1}{S_k}} = \frac{X_k/n}{S_k/n} \left(1 - O\left(\frac{1}{X_k}\right) + O\left(\frac{1}{S_k}\right) \right) = \frac{X_k/n}{S_k/n} + O\left(\frac{1}{S_k}\right),$$

since $S_k \geq X_k$. When we do the same computation for all the cases of Figure 3, we also get an error $O\left(\frac{1}{S_k}\right)$. Note that for the last three cases on the bottom right of Figure 3, the probability is already $O\left(\frac{1}{S_k}\right)$, so these cases do not contribute to $\phi_X(x, y, z)$. So we can bound (23) by

$$\frac{1}{\varepsilon_{k+1}^{3/4} \sqrt{n}} O\left(\frac{1}{S_k}\right) \stackrel{\text{Lem.1}}{=} O\left(\frac{1}{\varepsilon_{k+1}^{3/4} \sqrt{n}} \times \frac{1}{\varepsilon_k n}\right) \stackrel{\varepsilon_k \geq n^{-\frac{2}{5}-\frac{1}{100}}}{=} O\left(\frac{1}{\varepsilon_k n} \times \frac{1}{(n^{-\frac{2}{5}-\frac{1}{100}})^{3/4} \sqrt{n}}\right) = O\left(\frac{n^{-1/30}}{\varepsilon_k n}\right).$$

We move on to (24). We want to estimate the error when we do a linear approximation of ϕ_X near $(\mathcal{X}, \mathcal{Y}, \mathcal{Z})\left(\frac{k}{n}\right)$, so we will need to bound the second derivatives of ϕ_X near this point. More precisely, we write $(v_1, v_2, v_3) = \left(\frac{A_k}{n}, \frac{B_k}{n}, \frac{C_k}{n}\right)$. By the Taylor-Lagrange formula we can bound (24) by

$$\frac{1}{\varepsilon_{k+1}^{3/4} \sqrt{n}} \sum_{1 \leq i, j \leq 3} |v_i| \times |v_j| \times \max_{\substack{|u_1 - \mathcal{X}(k/n)| \leq |v_1| \\ |u_2 - \mathcal{Y}(k/n)| \leq |v_2| \\ |u_3 - \mathcal{Z}(k/n)| \leq |v_3|}} \left| \frac{\partial^2 \phi_X}{\partial x_i \partial x_j}(u_1, u_2, u_3) \right|. \quad (27)$$

By the assumptions of the proposition, we have the bounds:

$$|v_1| \leq 1000 \varepsilon_k^{3/4} \frac{\log n}{\sqrt{n}}, \quad |v_2| \leq 1000 \frac{\log n}{\sqrt{n}}, \quad |v_3| \leq 1000 \varepsilon_k^{1/2} \frac{\log n}{\sqrt{n}}. \quad (28)$$

On the other hand, we can compute the second order derivatives of ϕ_X , which are of the form $\frac{P(x,y,z)}{(x+y+z)^4}$ for some polynomial P . By Lemma 1, we know that u_1 , u_2 and u_3 are respectively $O(\varepsilon_k^2 n^{1/100})$, $O(\varepsilon_k n^{1/100})$ and $O(\varepsilon_k^{3/2} n^{1/100})$, and the sum $u_1 + u_2 + u_3$ is of order ε_k . Hence, we can consider the term with the highest order in the numerator. For example, we find

$$\frac{\partial^2 \phi_X}{\partial x_1^2} = \frac{12u_2^2 + 28u_2u_3 + 10u_3^2}{(u_1 + u_2 + u_3)^4},$$

and the highest order term in the numerator is $u_2^2 = O(\varepsilon_k^2 n^{1/50})$. On the other hand, the denominator is of order ε_k^4 , so we get

$$\frac{\partial^2 \phi_X}{\partial x_1^2}(u_1, u_2, u_3) = O(\varepsilon_k^{-2} n^{1/50}).$$

The bounds on $\frac{\partial^2 \phi_X}{\partial x_i \partial x_j}(u_1, u_2, u_3)$ that we obtain for all second-order partial derivatives are summarized in the following table:

$i \setminus j$	1	2	3
1	$O(\varepsilon_k^{-2} n^{1/50})$	$O(\varepsilon_k^{-3/2} n^{1/50})$	$O(\varepsilon_k^{-2} n^{1/50})$
2	$O(\varepsilon_k^{-3/2} n^{1/50})$	$O(\varepsilon_k^{-1} n^{1/50})$	$O(\varepsilon_k^{-3/2} n^{1/50})$
3	$O(\varepsilon_k^{-2} n^{1/50})$	$O(\varepsilon_k^{-3/2} n^{1/50})$	$O(\varepsilon_k^{-2} n^{1/50})$

Combining this with (28), we find that each term of (27) is

$$O\left(\frac{\varepsilon_k^{-1} n^{1/50} \log^2 n}{\varepsilon_{k+1}^{3/4} n \sqrt{n}}\right) = O\left(\frac{1}{\varepsilon_k n} \times \frac{n^{1/50} \log^2 n}{\varepsilon_{k+1}^{3/4} \sqrt{n}}\right) \stackrel{2\varepsilon_{k+1} \geq n^{-2/5-1/100}}{=} O\left(\frac{n^{-1/30}}{\varepsilon_k n}\right),$$

which bounds (24). Note that it was necessary to handle one by one the terms of (27) and not to bound everything crudely by $\|v\|^2 \times \|D^2 \phi_X\|$ (we would have obtained an additional factor ε_k^{-1} , which is too large).

Let us now bound (25). We first compute the gradient of $\nabla \phi_X$:

$$\nabla \phi_X(x, y, z) = \frac{1}{(x+y+z)^3} (-4y^2 - 9yz + xz - 3z^2, 4xy + 6xz + 2z^2, -x^2 - 2yz + 3xy + 3xz). \quad (29)$$

On the other hand, by (8), when $\varepsilon \rightarrow 0$, we have

$$\mathcal{X}(t_{\text{ext}} - \varepsilon) \sim 3\varepsilon^2, \quad \mathcal{Y}(t_{\text{ext}} - \varepsilon) \sim 4\varepsilon, \quad \mathcal{Z}(t_{\text{ext}} - \varepsilon) \sim 4\sqrt{3}\varepsilon^{3/2}.$$

Therefore, we can replace (x, y, z) in (29) by $(\mathcal{X}(t), \mathcal{Y}(t), \mathcal{Z}(t))$ and let $t \rightarrow t_{\text{ext}}$. We find that there are constants $K, \eta > 0$ such that, for any $0 < \varepsilon < \eta$, we have:

$$\begin{aligned} \left| \frac{\partial \phi_X}{\partial x} ((\mathcal{X}, \mathcal{Y}, \mathcal{Z})(t_{\text{ext}} - \varepsilon)) - \frac{1}{\varepsilon} \right| &\leq \frac{\delta}{\varepsilon}, \\ \left| \frac{\partial \phi_X}{\partial y} ((\mathcal{X}, \mathcal{Y}, \mathcal{Z})(t_{\text{ext}} - \varepsilon)) \right| &\leq K, \\ \left| \frac{\partial \phi_X}{\partial y} ((\mathcal{X}, \mathcal{Y}, \mathcal{Z})(t_{\text{ext}} - \varepsilon)) \right| &\leq \frac{K}{\varepsilon^{1/2}}. \end{aligned}$$

This is the value of η that we take in Proposition 4. We can now replace ε by $\varepsilon_k \in (0, \eta)$ and we obtain the following bound on (25):

$$\begin{aligned} & \frac{\delta}{\varepsilon_k n} \times \frac{1}{\varepsilon_{k+1}^{3/4} \sqrt{n}} |A_k| + \frac{K}{n} \times \frac{1}{\varepsilon_{k+1}^{3/4} \sqrt{n}} |B_k| + \frac{K}{\varepsilon_k^{1/2} n} \times \frac{1}{\varepsilon_{k+1}^{3/4} \sqrt{n}} |C_k| \\ &= \frac{\delta}{\varepsilon_k n} |\tilde{A}_k| + \frac{K}{\varepsilon_k n} \varepsilon_k^{1/4} |\tilde{B}_k| + \frac{K}{\varepsilon_k n} \varepsilon_k^{1/4} |\tilde{C}_k| \\ &\leq \frac{\delta}{\varepsilon_k n} |\tilde{A}_k| + \frac{1000K}{\varepsilon_k n} \varepsilon_k^{1/4} \log n. \end{aligned}$$

We finally treat the term (26). We recall that \mathcal{X} solves the equation $\mathcal{X}' = \phi_X(\mathcal{X}, \mathcal{Y}, \mathcal{Z})$, so this is just a linear approximation, so we will need to bound the second derivative \mathcal{X}'' . More precisely, (26) is bounded by

$$\frac{n}{\varepsilon_{k+1}^{3/4} \sqrt{n}} \times \left(\frac{1}{n}\right)^2 \times \max_{\left[\frac{k}{n}, \frac{k+1}{n}\right]} |\mathcal{X}''|. \quad (30)$$

Moreover, by differentiating $\mathcal{X}' = \phi_X(\mathcal{X}, \mathcal{Y}, \mathcal{Z})$, we have

$$\begin{aligned} \mathcal{X}''(t) &= \left(\phi_X \frac{\partial \phi_X}{\partial x} + \phi_Y \frac{\partial \phi_X}{\partial y} + \phi_Z \frac{\partial \phi_X}{\partial z} \right) (\mathcal{X}(t), \mathcal{Y}(t), \mathcal{Z}(t)) \\ &= \frac{\mathcal{Z}(t)}{\mathcal{I}(t)^4} (\mathcal{X}(t)^2 - 2\mathcal{X}(t)\mathcal{Y}(t) + 8\mathcal{Y}(t)\mathcal{Z}(t) + 11\mathcal{Z}(t)^2). \end{aligned}$$

This is a continuous function of t on $[0, t_{\text{ext}})$. Moreover, by (8), we have

$$\mathcal{X}(t_{\text{ext}} - \varepsilon) \sim_{\varepsilon \rightarrow 0} \frac{4\sqrt{3}\varepsilon^{3/2}}{(4\varepsilon)^4} \times 8 \times 4\varepsilon \times 4\sqrt{3}\varepsilon^{3/2} = 6,$$

so \mathcal{X}'' is bounded by a constant K . Plugging this into (30), we can bound (26) by $\frac{K}{\varepsilon_{k+1}^{3/4} n^{3/2}} = O\left(\frac{n^{-1/30}}{\varepsilon_k n}\right)$.

WE NOW MOVE ON TO THE ESTIMATES (16) AND (17). Since the proof is similar, we will not do it in full details and only stress the differences with the proof of (15). The decomposition of the error into five terms is the same with the following modifications:

- the first term (22) becomes $\left| \left(\frac{\varepsilon_k^{1/2}}{\varepsilon_{k+1}^{1/2}} - 1 \right) \tilde{C}_k - \frac{1}{2} \frac{1}{\varepsilon_k n} \tilde{C}_k \right|$ for \tilde{C} , and disappears completely for \tilde{B} ;
- in the terms (23), (24), (25) and (26), the factors $\frac{1}{\varepsilon_{k+1}^{3/4} \sqrt{n}}$ become $\frac{1}{\sqrt{n}}$ for \tilde{B} and $\frac{1}{\varepsilon_{k+1}^{1/2} \sqrt{n}}$ for \tilde{C} ;
- in the fourth term (25), the drift $\frac{1}{\varepsilon_k n} \tilde{A}_k$ becomes 0 for \tilde{B} and $\frac{3\sqrt{3}}{2} \tilde{B}_k - \frac{3}{2} \tilde{C}_k$ of \tilde{C} .

The first and second term can then be bounded by $O\left(\frac{n^{-1/30}}{\varepsilon_k n}\right)$ in the exact same way as for \tilde{A} (this bound actually becomes cruder for (23), since now the factor $\varepsilon_{k+1}^{3/4}$ in the denominator disappears or become larger).

The bound on the fifth term (26) is also very similar: we now have

$$\begin{aligned} \mathcal{Y}''(t) &= -2 \frac{\mathcal{Z}}{\mathcal{I}^4} (\mathcal{X}\mathcal{Y} + 4\mathcal{Y}^2 + 8\mathcal{X}\mathcal{Z} + 21\mathcal{Y}\mathcal{Z} + 20\mathcal{Z}^2)(t) = O\left((t_{\text{ext}} - t)^{-1/2}\right) \\ \mathcal{X}''(t) &= 3 \frac{\mathcal{Z}}{\mathcal{I}^4} (\mathcal{X}^2 + 4\mathcal{X}\mathcal{Y} + 4\mathcal{Y}^2 + 8\mathcal{X}\mathcal{Z} + 14\mathcal{Y}\mathcal{Z} + 11\mathcal{Z}^2)(t) = O\left((t_{\text{ext}} - t)^{-1/2}\right). \end{aligned}$$

Therefore, the analog of (26) for \tilde{B} (resp. \tilde{C}) is $O\left(\frac{1}{\sqrt{n}} \times n \times \left(\frac{1}{n}\right)^2 \times \varepsilon_k^{-1/2}\right) = O\left(\frac{\varepsilon_k^{-1/2}}{n^{3/2}}\right)$ (resp. $O\left(\frac{\varepsilon_k^{-1}}{n^{3/2}}\right)$). In both cases, this is $O\left(\frac{n^{-1/30}}{\varepsilon_k n}\right)$.

The analog of the third term (24) is still very similar, but requires to be more careful. Indeed (27) becomes respectively

$$\frac{1}{\varepsilon_{k+1}^{1/2} \sqrt{n}} \sum_{1 \leq i, j \leq 3} |v_i| \times |v_j| \times \max_{\substack{|u_1 - \mathcal{X}(k/n)| \leq |v_1| \\ |u_2 - \mathcal{Y}(k/n)| \leq |v_2| \\ |u_3 - \mathcal{Z}(k/n)| \leq |v_3|}} \left| \frac{\partial^2 \phi_Y}{\partial x_i \partial x_j}(u_1, u_2, u_3) \right|. \quad (31)$$

and

$$\frac{1}{\sqrt{n}} \sum_{1 \leq i, j \leq 3} |v_i| \times |v_j| \times \max_{\substack{|u_1 - \mathcal{X}(k/n)| \leq |v_1| \\ |u_2 - \mathcal{Y}(k/n)| \leq |v_2| \\ |u_3 - \mathcal{Z}(k/n)| \leq |v_3|}} \left| \frac{\partial^2 \phi_Z}{\partial x_i \partial x_j}(u_1, u_2, u_3) \right|. \quad (32)$$

for \tilde{B} and \tilde{C} . Moreover, when we compute the second order partial derivatives $\frac{\partial^2 \phi_Y}{\partial x_i \partial x_j}$ and $\frac{\partial^2 \phi_Z}{\partial x_i \partial x_j}$, we get respectively the following tables:

$i \setminus j$	1	2	3
1	$O(\varepsilon_k^{-2} n^{1/50})$	$O(\varepsilon_k^{-2} n^{1/50})$	$O(\varepsilon_k^{-2} n^{1/50})$
2	$O(\varepsilon_k^{-2} n^{1/50})$	$O(\varepsilon_k^{-3/2} n^{1/50})$	$O(\varepsilon_k^{-2} n^{1/50})$
3	$O(\varepsilon_k^{-2} n^{1/50})$	$O(\varepsilon_k^{-2} n^{1/50})$	$O(\varepsilon_k^{-3/2} n^{1/50})$

$i \setminus j$	1	2	3
1	$O(\varepsilon_k^{-3/2} n^{1/50})$	$O(\varepsilon_k^{-3/2} n^{1/50})$	$O(\varepsilon_k^{-2} n^{1/50})$
2	$O(\varepsilon_k^{-3/2} n^{1/50})$	$O(\varepsilon_k^{-3/2} n^{1/50})$	$O(\varepsilon_k^{-2} n^{1/50})$
3	$O(\varepsilon_k^{-2} n^{1/50})$	$O(\varepsilon_k^{-2} n^{1/50})$	$O(\varepsilon_k^{-2} n^{1/50})$

In both cases, using (28), we find that each term of (31) or (32) is

$$O\left(\frac{\varepsilon_k^{-3/2} n^{1/50} \log^2 n}{\varepsilon_{k+1}^{1/2} n^{3/2}}\right) = O\left(\frac{1}{\varepsilon_k n} \times \frac{n^{1/50} \log^2 n}{\varepsilon_k \sqrt{n}}\right) \underset{\varepsilon_k \geq n^{-\frac{2}{5} - \frac{1}{100}}}{=} O\left(\frac{n^{-1/30}}{\varepsilon_k n}\right).$$

Finally, to handle the analog of the fourth term (25), we just need to compute the gradients of ϕ_Y and ϕ_Z :

$$\nabla \phi_Y(x, y, z) = \frac{1}{(x + y + z)^3} (2xy + 6y^2 + 6yz - 8z^2, -2x^2 - 6xy - 6xz + 4yz + 12z^2, 8xz + 4y^2 + 12yz),$$

$$\nabla \phi_Z(x, y, z) = \frac{1}{(x + y + z)^3} (3xz + 9yz + 15z^2, 6yz + 12z^2, -3x^2 - 9xy - 15xz - 6y^2 - 12yz).$$

As in the first case, we can now replace (x, y, z) by $(\mathcal{X}(t), \mathcal{Y}(t), \mathcal{Z}(t))$, use (8) and identify the highest order terms in $t_{\text{ext}} - t$. We find that there is a constant K such that for all $0 \leq t < t_{\text{ext}}$:

$$\begin{aligned} \left| \frac{\partial \phi_Y}{\partial x}(\mathcal{X}(t), \mathcal{Y}(t), \mathcal{Z}(t)) \right| &\leq \frac{K}{t_{\text{ext}} - t}, \\ \left| \frac{\partial \phi_Y}{\partial y}(\mathcal{X}(t), \mathcal{Y}(t), \mathcal{Z}(t)) \right| &\leq \frac{K}{(t_{\text{ext}} - t)^{1/2}}, \\ \left| \frac{\partial \phi_Y}{\partial z}(\mathcal{X}(t), \mathcal{Y}(t), \mathcal{Z}(t)) \right| &\leq \frac{K}{t_{\text{ext}} - t}, \\ \left| \frac{\partial \phi_Z}{\partial x}(\mathcal{X}(t), \mathcal{Y}(t), \mathcal{Z}(t)) \right| &\leq \frac{K}{(t_{\text{ext}} - t)^{1/2}}. \end{aligned}$$

Moreover, there is $\eta > 0$ (depending on δ) such that, if $t_{\text{ext}} - \eta \leq t < t_{\text{ext}}$, then

$$\begin{aligned} \left| \frac{\partial \phi_Z}{\partial y} (\mathcal{X}(t), \mathcal{Y}(t), \mathcal{Z}(t)) - \frac{3\sqrt{3}}{2} \frac{1}{(t_{\text{ext}} - t)^{1/2}} \right| &\leq \frac{\delta}{(t_{\text{ext}} - t)^{1/2}}, \\ \left| \frac{\partial \phi_Z}{\partial y} (\mathcal{X}(t), \mathcal{Y}(t), \mathcal{Z}(t)) + \frac{3}{2} \frac{1}{t_{\text{ext}} - t} \right| &\leq \frac{\delta}{t_{\text{ext}} - t}. \end{aligned}$$

From here, taking $t = \frac{k}{n}$ and replacing (A_k, B_k, C_k) by $(\varepsilon_k^{3/4} \sqrt{n} \tilde{A}_k, \sqrt{n} \tilde{B}_k, \varepsilon_k^{1/2} \sqrt{n} \tilde{C}_k)$, we easily obtain the claimed bound on (25). \square

Proof of Proposition 5. Just like in the proof of Proposition 4, we first introduce the following functions (again with the notation $s = x + y + z$):

$$\begin{aligned} \psi_X(x, y, z) &= 4\frac{x}{s} + 4\frac{xy}{s^2} + \frac{yz}{s^2} + \frac{y^2z}{s^3} + 9\frac{x^2z}{s^3} + 2\frac{xyz}{s^3} + 2\frac{xz^2}{s^3} + \frac{z^3}{s^3}, \\ \psi_Y(x, y, z) &= \frac{xy}{s^2} + 4\frac{y^2}{s^2} + 4\frac{y^2z}{s^3} + 2\frac{xyz}{s^3} + 2\frac{xz^2}{s^3} + 4\frac{z^3}{s^3}, \\ \psi_Z(x, y, z) &= \frac{yz}{s^2} + \frac{y^2z}{s^3} + 8\frac{yz^2}{s^3} + \frac{x^2z}{s^3} + 2\frac{xyz}{s^3} + 8\frac{xz^2}{s^3} + 9\frac{z^3}{s^3}. \end{aligned}$$

These functions are respectively the fluid limit approximations of $\mathbb{E}[(\Delta X_k)^2 | \mathcal{F}_k]$, $\mathbb{E}[(\Delta Y_k)^2 | \mathcal{F}_k]$ and $\mathbb{E}[(\Delta Z_k)^2 | \mathcal{F}_k]$ and can be computed from the transitions given in Figure 3 as before.

VARIANCE OF \tilde{A} . Let us start by establishing (18). We first note that, since adding a function of A_k does not change the conditional variance on \mathcal{F}_k , we have

$$\text{Var}(\Delta \tilde{A}_k | \mathcal{F}_k) = \text{Var}\left(\Delta \tilde{A}_k + \left(\frac{1}{\varepsilon_k^{3/4} \sqrt{n}} - \frac{1}{\varepsilon_{k+1}^{3/4} \sqrt{n}}\right) A_k | \mathcal{F}_k\right) = \frac{1}{\varepsilon_{k+1}^{3/2} n} \text{Var}(\Delta A_k | \mathcal{F}_k) = \frac{1}{\varepsilon_{k+1}^{3/2} n} \text{Var}(\Delta X_k | \mathcal{F}_k).$$

Therefore, we can write

$$\text{Var}(\Delta \tilde{A}_k | \mathcal{F}_k) = \frac{1}{\varepsilon_{k+1}^{3/2} n} \mathbb{E}[(\Delta X_k)^2 | \mathcal{F}_k] - \frac{1}{\varepsilon_{k+1}^{3/2} n} \mathbb{E}[\Delta X_k | \mathcal{F}_k]^2,$$

so

$$\begin{aligned} \left| \text{Var}(\Delta \tilde{A}_k | \mathcal{F}_k) - \frac{2\sqrt{3}}{\varepsilon_k n} \right| &\leq \frac{1}{\varepsilon_{k+1}^{3/2} n} \mathbb{E}[\Delta X_k | \mathcal{F}_k]^2 \\ &+ \frac{1}{\varepsilon_{k+1}^{3/2} n} \left| \mathbb{E}[(\Delta X_k)^2 | \mathcal{F}_k] - \psi_X\left(\frac{X_k}{n}, \frac{Y_k}{n}, \frac{Z_k}{n}\right) \right| \\ &+ \frac{1}{\varepsilon_{k+1}^{3/2} n} \left| \psi_X\left(\frac{X_k}{n}, \frac{Y_k}{n}, \frac{Z_k}{n}\right) - \psi_X\left(\mathcal{X}\left(\frac{k}{n}\right), \mathcal{Y}\left(\frac{k}{n}\right), \mathcal{Z}\left(\frac{k}{n}\right)\right) \right| \\ &+ \frac{1}{\varepsilon_{k+1}^{3/2} n} \left| \psi_X\left(\mathcal{X}\left(\frac{k}{n}\right), \mathcal{Y}\left(\frac{k}{n}\right), \mathcal{Z}\left(\frac{k}{n}\right)\right) - (2\sqrt{3}\sqrt{\varepsilon_k}) \right| \\ &+ \frac{1}{\varepsilon_{k+1}^{3/2} n} \left| (2\sqrt{3}\sqrt{\varepsilon_k}) - (2\sqrt{3}\sqrt{\varepsilon_{k+1}}) \right| \end{aligned} \quad (33)$$

The first term can be bounded by using Proposition 4 and $\varepsilon_k \geq n^{-2/5-1/100}$. Moreover, by the exact same argument as for term (23) in the proof of Proposition 4, the second term is

$$O\left(\frac{1}{\varepsilon_{k+1}^{3/2} n} \times \frac{1}{S_k}\right) \stackrel{\text{Lem.1}}{=} O\left(\frac{1}{\varepsilon_k^{5/2} n^2}\right) \stackrel{\varepsilon_k \geq n^{-\frac{2}{5}-\frac{1}{100}}}{=} O\left(\frac{n^{-1/30}}{\varepsilon_k n}\right).$$

We now bound the third term of (33). If we write $(v_1, v_2, v_3) = \left(\frac{A_k}{n}, \frac{B_k}{n}, \frac{C_k}{n}\right)$, this is bounded by

$$\frac{1}{\varepsilon_{k+1}^{3/2} \sqrt{n}} \sum_{i=1}^3 |v_i| \times \max_{\substack{|x-\mathcal{X}(k/n)| \leq |v_1| \\ |y-\mathcal{Y}(k/n)| \leq |v_2| \\ |z-\mathcal{Z}(k/n)| \leq |v_3|}} \left| \frac{\partial \psi_X}{\partial x}(x, y, z) \right|. \quad (34)$$

Just like for ϕ_X in the proof of Proposition 4, we can compute the gradient of ψ_X : the partial derivatives are of the form $\frac{P(x,y,z)}{(x+y+z)^4}$, where P is a homogeneous polynomial of degree 3. By using Lemma 1, just like in (27), we have that x , y and z are respectively $O(\varepsilon_k^2 n^{1/100})$, $O(\varepsilon_k n^{1/100})$ and $O(\varepsilon_k^{3/2} n^{1/100})$ and that the sum $x+y+z$ is of order ε_k . Therefore, by considering the higher order terms in the polynomial $P(x, y, z)$, we obtain the following estimates:

$$\frac{\partial \psi_X}{\partial x}(x, y, z) = O(\varepsilon_k^{-1} n^{3/100}), \quad \frac{\partial \psi_X}{\partial y}(x, y, z) = O(\varepsilon_k^{-1/2} n^{3/100}), \quad \frac{\partial \psi_X}{\partial z}(x, y, z) = O(\varepsilon_k^{-1} n^{3/100}).$$

Combining this with (34), we get that the third term of (33) is

$$O\left(\frac{n^{3/100} \log n}{\varepsilon_k^2 n^{3/2}}\right).$$

using $\varepsilon_k \geq n^{-2/5-1/100}$, this is $O\left(\frac{n^{-1/30}}{\varepsilon_k n}\right)$.

We now bound the fourth term of (33). For this, we use again (8). In particular, when we write down $\psi_X(\mathcal{X}, \mathcal{Y}, \mathcal{Z})(t_{\text{ext}} - \varepsilon)$, the highest order terms in ε are $\frac{\mathcal{Y}\mathcal{Z}}{\mathcal{X}^2} \sim \sqrt{3}\sqrt{\varepsilon}$ and $\frac{\mathcal{Y}^2\mathcal{Z}}{\mathcal{X}^3} \sim \sqrt{3}\sqrt{\varepsilon}$, so we have

$$\psi_X(\mathcal{X}, \mathcal{Y}, \mathcal{Z})(t_{\text{ext}} - \varepsilon) \sim_{\varepsilon \rightarrow 0} 2\sqrt{3}\sqrt{\varepsilon}.$$

In particular, taking $\varepsilon = \varepsilon_k$, if we choose η small enough the fourth term of (33) is bounded by $\frac{\delta}{\varepsilon_k n}$. Finally the fifth term is also smaller than $\frac{\delta}{\varepsilon_k n}$ if ε_k is small enough. Gathering-up the pieces we have established (18).

The bound (19) follows from the same proof by noticing that the only term of (33) which makes the errors $\frac{\tilde{A}^2}{\varepsilon_k^{1/2} n}$, $\frac{\tilde{B}^2}{n}$, $\frac{\tilde{C}^2}{n}$ appear is the $-\mathbb{E}[\Delta X_k | \mathcal{F}_k]^2$, which is negative.

VARIANCE OF \tilde{B} . The bound (20) is immediate: for the same reason as with \tilde{A} , we have

$$\text{Var}\left(\Delta \tilde{B}_k | \mathcal{F}_k\right) = \frac{1}{n} \text{Var}\left(\Delta Y_k | \mathcal{F}_k\right) \leq \frac{9}{n},$$

since $|\Delta Y_k|$ is bounded by 3.

VARIANCE OF \tilde{C} . Finally, we prove (21): as before, we can write

$$\begin{aligned} \text{Var}\left(\Delta \tilde{C}_k | \mathcal{F}_k\right) &= \frac{1}{\varepsilon_{k+1} n} \mathbb{E}\left[(\Delta Z_k)^2 | \mathcal{F}_k\right] - \mathbb{E}\left[\Delta Z_k | \mathcal{F}_k\right]^2 \\ &\leq \frac{1}{\varepsilon_{k+1} n} \left| \mathbb{E}\left[(\Delta Z_k)^2 | \mathcal{F}_k\right] - \Psi_Z\left(\frac{X_k}{n}, \frac{Y_k}{n}, \frac{Z_k}{n}\right) \right| + \frac{1}{\varepsilon_{k+1} n} \Psi_Z\left(\frac{X_k}{n}, \frac{Y_k}{n}, \frac{Z_k}{n}\right). \end{aligned}$$

By the exact same argument as for \tilde{A} , the first term is

$$O\left(\frac{1}{\varepsilon_{k+1} n} \frac{1}{S_k}\right) = O\left(\frac{1}{\varepsilon_k^2 n^2}\right) = O\left(\frac{\varepsilon_k^{1/2}}{\varepsilon_k n}\right),$$

where the first equality comes from Lemma (1) and the second from $\varepsilon_k \geq n^{-2/5-1/100}$. On the other hand, noticing that every term in $\psi_Z(x, y, z)$ has a factor $\frac{z}{s}$, we can write

$$\psi_Z\left(\frac{X_k}{n}, \frac{Y_k}{n}, \frac{Z_k}{n}\right) = O\left(\frac{Z_k}{S_k}\right) = O\left(\frac{\varepsilon_k^{3/2} n}{\varepsilon_k n}\right) = O\left(\varepsilon_k^{1/2}\right),$$

where the second inequality comes from Lemma 1. This proves (21). □

4.2 Rough behaviour of \tilde{A} , \tilde{B} and \tilde{C}

In this section we will use our drift and variance estimates to control \tilde{A} , \tilde{B} , \tilde{C} . Recall notation from (14) and (13). We shall get a rather rough control on \tilde{A} , \tilde{B} and \tilde{C} (Proposition 6) and later refine the one on \tilde{A} . In the rest of this subsection, on top of the constant $K > 0$ given by Propositions 4 and 5, we fix

$$\delta = \frac{1}{100}$$

for definiteness and let $0 < \eta \equiv \eta(\delta) < 1/2$ so that we can apply the above propositions. In particular, the value of η does not depend on n , nor on the coming $\epsilon > 0$ and *its value may be decreased for convenience* by keeping the same δ . In the coming pages $\text{Cst} > 0$ is a constant (which may depend on the constant K or the now-fixed $\delta = \frac{1}{100}$) and that may increase from line to line, but whose value does not depend upon n (provided it is large enough), nor η , nor on the forthcoming ϵ . On the contrary K_ϵ is a constant that depends upon ϵ but also upon η in an implicit way.

The value η shall give our “starting scale” $k_0 = \lfloor (t_{\text{ext}} - \eta)n \rfloor$ which is such that $\varepsilon_{k_0} = \eta$ and we shall then look at times $k_0 \leq k \leq \tilde{\theta}$. We start by controlling the fluctuations at k_0 .

Lemma 2 (Fluctuations in the bulk). *For all $\epsilon > 0$ there exists $K_\epsilon > 0$ so that for all n large enough, with probability at least $1 - \epsilon$ we have*

$$\max(|\tilde{A}_{k_0}|, |\tilde{B}_{k_0}|, |\tilde{C}_{k_0}|) < K_\epsilon \text{ and } \tilde{\theta} > k_0. \quad (35)$$

Proof. Classical results entail that on top of the law of large numbers for the process $n^{-1} \cdot (X^n, Y^n, Z^n)$ given in Proposition 2, we have a functional central limit theorem for their fluctuations, as long as we stay in the bulk. More precisely, for $0 \leq t \leq (t_{\text{ext}} - \eta)$, the solution given by the differential equation (5) is bounded away from 0, i.e.

$$\inf\{\min(\mathcal{X}(t), \mathcal{Y}(t), \mathcal{Z}(t)) : 0 \leq t \leq (t_{\text{ext}} - \eta)\} > 0, \quad (36)$$

and thanks to our hypothesis (3), the initial fluctuations A_0 , B_0 and C_0 are bounded so that $(A_0, B_0, C_0)/\sqrt{n}$ converges to $(0, 0, 0)$ ⁶. Therefore, we can apply [7, Theorem 2.3 p 458], which implies that

$$\left(\left(\frac{A_{\lfloor tn \rfloor}}{\sqrt{n}}, \frac{B_{\lfloor tn \rfloor}}{\sqrt{n}}, \frac{C_{\lfloor tn \rfloor}}{\sqrt{n}} \right) : 0 \leq t \leq t_{\text{ext}} - \eta \right)$$

converges as n goes to infinity weakly to a continuous random processes driven by a nice stochastic differential equation. Furthermore [7, Theorem 2.3 p 458] entails that the terminal value

$$\left(\frac{A_{\lfloor (t_{\text{ext}} - \eta)n \rfloor}}{\sqrt{n}}, \frac{B_{\lfloor (t_{\text{ext}} - \eta)n \rfloor}}{\sqrt{n}}, \frac{C_{\lfloor (t_{\text{ext}} - \eta)n \rfloor}}{\sqrt{n}} \right)$$

converges towards a Gaussian law whose covariance depends on η only. Given (36), this implies that w.h.p. we have $X_k > 0$ for all $0 \leq k \leq (t_{\text{ext}} - \eta)n$ (in other words $\theta > (t_{\text{ext}} - \eta)n$) and that $|\tilde{A}_{\lfloor (t_{\text{ext}} - \eta)n \rfloor}|, |\tilde{B}_{\lfloor (t_{\text{ext}} - \eta)n \rfloor}|, |\tilde{C}_{\lfloor (t_{\text{ext}} - \eta)n \rfloor}|$ are tight. The statement of the lemma follows. □

⁶We could have allowed $o(\sqrt{n})$ fluctuations, but not $o(n)$ as in Theorem 1.

After this initial control, we shall provide a rough upper bound on the fluctuation processes.

Proposition 6 (Rough upper bounds). *For all $\epsilon > 0$, there exists a constant $K_\epsilon > 0$ such that for n large enough, with probability at least $1 - \epsilon$ we have*

$$\max_{k_0 \leq k < \tilde{\theta}} \left\{ \frac{|\tilde{A}_k|}{|\log(\epsilon_k)|^{3/4}}, |\tilde{B}_k|, |\tilde{C}_k| \right\} \leq K_\epsilon. \quad (37)$$

Remark (The truth). The proof of the proposition shows that we can replace the $3/4$ exponent by $1/2 + \delta$ for all $\delta > 0$. We anyway expect an “iterated logarithm” behavior for \tilde{A} so that we could replace $|\log(\epsilon_k)|$ by $|\log \log(\epsilon_k)|$. In the same vein, a little more effort would yield that \tilde{B} and \tilde{C} “converge”⁷ but our estimates will be largely sufficient for our purposes.

Proof. In light of the form of the drift of \tilde{C} obtained in Equation (17), we will rather consider the process $\tilde{E}_k = \tilde{C}_k - \frac{3\sqrt{3}}{2}\tilde{B}_k$ instead of \tilde{C} , but notice we can control $|\tilde{C}_k| \leq \frac{3\sqrt{3}}{2}|\tilde{B}_k| + |\tilde{E}_k|$ using the processes \tilde{B} and \tilde{E} so that it is sufficient to prove the proposition after replacing \tilde{C} by \tilde{E} . Introduce L the first time at which one of the those three processes becomes large, i.e.

$$L = \tilde{\theta} \wedge \min \left\{ k \geq k_0 : \max \left(\frac{|\tilde{A}_k|}{|\log \epsilon_k|^{3/4}}, |\tilde{B}_k|, |\tilde{E}_k| \right) > K_\epsilon \right\}.$$

We call the region defined by the above inequalities on $(\tilde{A}, \tilde{B}, \tilde{C})$ the *good region* for the processes and evaluate separately the probability that we exit this region (i.e. that $L < \tilde{\theta}$) via one of the three processes \tilde{A}, \tilde{B} or \tilde{E} . By definition (14) of $\tilde{\theta}$ and since we will always take n large enough to have

$$K_\epsilon(1 + \log_2^{3/4}(n)) < \log(n),$$

as long as $k_0 \leq k < L$, we can apply the estimates obtained in Propositions 4 and 5. Specifically, we will decompose the processes \tilde{A}, \tilde{B} and \tilde{E} into their predictable and martingale parts and use Doob’s maximal inequality and L^2 estimates to control the martingales.

LET US START WITH \tilde{B} . We write for $k_0 \leq k \leq L$,

$$\tilde{B}_k = \tilde{B}_{k_0} + \sum_{\ell=k_0}^{k-1} \mathbb{E} \left[\Delta \tilde{B}_\ell | \mathcal{F}_\ell \right] + M_k^B$$

where $(M_{k \wedge L}^B)_{k \geq k_0}$ is an (\mathcal{F}_k) -martingale which starts from 0 at time k_0 . We first evaluate the drift/predictable part. To ease the calculation and readability, we will deliberately drop the integer-part notation $\lfloor \cdot \rfloor$ and introduce scales. Recall that the value of $\eta = \epsilon_{k_0}$ has been fixed above, but we may decrease it for convenience as long as it is independent of n and ϵ . We start from $k_0 = (t_{\text{ext}} - \eta)n$ and we let

$$k_j = (t_{\text{ext}} - \eta 2^{-j})n,$$

⁷To be precise, and stressing the dependence in n , the processes $(\tilde{B}_{\lfloor tn \rfloor \wedge \theta^n}^n : t \in [0, t_{\text{ext}}])$ converge in law for the $\|\cdot\|_\infty$ distance towards a limiting process $(\mathcal{B}_t : t \in [0, t_{\text{ext}}])$ which is continuous and in particular continuous at t_{ext} . Similarly $(\tilde{C}_{\lfloor tn \rfloor \wedge \theta^n}^n : t \in [0, t_{\text{ext}}]) \rightarrow (\mathcal{C}_t : t \in [0, t_{\text{ext}}])$ for a random continuous process and furthermore $\mathcal{C}_{t_{\text{ext}}} = \frac{3\sqrt{3}}{2}\mathcal{B}_{t_{\text{ext}}}$.

for $0 \leq j \leq (\frac{2}{5} + \frac{1}{100}) \log_2(n)$. In particular we have $j + |\log_2 \eta| \leq |\log_2 \epsilon_k| \leq (j+1) + |\log_2 \eta|$ for all $k_j \leq k \leq k_{j+1}$. With this notation, and using our estimate (16), we know that if $k \geq k_0$ we have

$$\begin{aligned}
\mathbb{1}_{k \leq L} \left| \sum_{\ell=k_0}^{k-1} \mathbb{E} [\Delta \tilde{B}_\ell | \mathcal{F}_\ell] \right| &\stackrel{(16)}{\leq} \text{Cst} \sum_{\ell=k_0}^{\infty} \mathbb{1}_{\ell < L} \left(\frac{\max(|\tilde{A}_\ell|, |\tilde{B}_\ell|, |\tilde{C}_\ell|)}{\sqrt{\epsilon_\ell n}} + \frac{1}{\epsilon_\ell n} n^{-1/30} \right) \\
&\leq_{\text{good region}} \text{Cst} \cdot K_\epsilon \cdot \sum_{\ell=k_0}^{\infty} \mathbb{1}_{\ell < L} \left(\frac{|\log \epsilon_\ell|^{3/4} + 1}{\sqrt{\epsilon_\ell n}} + \frac{1}{\epsilon_\ell n} n^{-1/30} \right) \\
&\leq_{\text{scales}} \text{Cst} \cdot K_\epsilon \cdot \sum_{j=1}^{\log_2(n)} \sum_{\ell=k_{j-1}}^{k_j-1} \left(\frac{(j + |\log_2 \eta|)^{3/4}}{\sqrt{\eta 2^{-j} n}} + \frac{1}{\eta 2^{-j} n} n^{-1/30} \right) \\
&\leq \text{Cst} \cdot K_\epsilon \cdot \sum_{j=1}^{\log_2(n)} \left(\sqrt{\eta} \frac{(j + |\log_2 \eta|)^{3/4}}{2^{j/2}} + n^{-1/30} \right) \\
&\leq \text{Cst} \cdot K_\epsilon \cdot |\log \eta| \cdot \sum_{j=1}^{\log_2(n)} \left(\sqrt{\eta} \frac{j+1}{2^{j/2}} + n^{-1/30} \right) \\
&\leq K_\epsilon \cdot (\text{Cst} \cdot \sqrt{\eta} |\log \eta|),
\end{aligned}$$

for n large enough where $\text{Cst} > 0$ is a constant that may vary from line to line but that does not depend on n , nor on ϵ nor on η as long as it is small. In particular, we may decrease the value of η so that the parenthesis in the last display is smaller than $1/4$ say. We obtain that the sum of the absolute values of the expected conditional drifts of \tilde{B} between k_0 and L is bounded by $K_\epsilon/4$ (deterministically).

We deduce that the event $\{L < \tilde{\theta} \text{ and } |B_{k_0}| \leq K_\epsilon/4 \text{ and } |\tilde{B}_L| > K_\epsilon\}$ is included in the event $\{L < \tilde{\theta} \text{ and } |M_L^B| > K_\epsilon/2\}$ so that in particular we can write

$$\begin{aligned}
\mathbb{P}(L < \tilde{\theta} \text{ and we exit the region by } \tilde{B}) &\leq \mathbb{P}(|\tilde{B}_{k_0}| > K_\epsilon/4) + \mathbb{P}(L < \tilde{\theta} \text{ and } |M_L^B| > K_\epsilon/2) \\
&\leq \mathbb{P}(|\tilde{B}_{k_0}| > K_\epsilon/4) + \mathbb{P}\left(\sup_{k_0 \leq k < L} |M_k^B| > K_\epsilon/2\right) \\
&\leq_{\text{Doob}} \mathbb{P}(|\tilde{B}_{k_0}| > K_\epsilon/4) + 4 \frac{\mathbb{E}[(M_L^B)^2]}{K_\epsilon^2/4}.
\end{aligned}$$

Up to increasing K_ϵ we can bound the first term by ϵ using Lemma 2. To bound the second term, we use our variance estimate (20) which gives in the good region

$$\mathbb{E}[(\Delta M_k^B)^2 | \mathcal{F}_k, k \leq L] = \text{Var}(\Delta M_k^B | \mathcal{F}_k, k \leq L) = \text{Var}(\Delta \tilde{B}_k | \mathcal{F}_k, k \leq L) \stackrel{(20)}{\leq} \frac{K}{n}.$$

By the orthogonality of martingale increments in L^2 we deduce that

$$\mathbb{E}[(M_L^B)^2] = \sum_{k=k_0}^{\infty} \mathbb{E}[(\Delta M_k^B)^2 \mathbb{1}_{k \leq L} | \mathcal{F}_k] \leq \frac{K(t_{\text{ext}} n - k_0)}{n} = K\eta.$$

Hence we obtain

$$\mathbb{P}(L < \tilde{\theta} \text{ and we exit the good region by } \tilde{B}) \leq \epsilon + 16 \frac{\mathbb{E}[(M_L^B)^2]}{K_\epsilon^2} \leq \epsilon + \frac{16K\eta}{K_\epsilon^2}.$$

If K_ϵ is large enough, the second term is also less than ϵ so that the probability in the left-hand side is small. Conclusion: it is unlikely that we exit first the good region because of the process \tilde{B} .

CASE OF \tilde{E} . The proof is similar, but we shall use more precisely the form of the conditional expected drifts. As before, we write

$$\tilde{E}_k = \tilde{E}_{k_0} + \sum_{\ell=k_0}^{k-1} \mathbb{E} \left[\Delta \tilde{E}_\ell | \mathcal{F}_\ell \right] + M_k^E$$

where $(M_{k \wedge L}^E)_{k \geq k_0}$ is an (\mathcal{F}_k) -martingale which starts from 0 at time k_0 . We will bound $\mathbb{P}(L < \tilde{\theta} \text{ and } \tilde{E}_L > K_\epsilon)$ and the case $\tilde{E}_L < -K_\epsilon$ will be treated similarly. Let us introduce L_E^- , the last time before L where \tilde{E} is smaller than $K_\epsilon/2$. In particular on the event $\{L < \tilde{\theta} \text{ and } \tilde{E}_L > K_\epsilon\}$, for $L_E^- < k \leq L$ the process \tilde{E} is larger than $K_\epsilon/2$ and its conditional expected drift therefore satisfies

$$\begin{aligned} \left| \mathbb{E} \left[\Delta \tilde{E}_k | \mathcal{F}_k \right] - \frac{1}{\epsilon_k n} \underbrace{\tilde{E}_k}_{\geq K_\epsilon/2} \right| &\stackrel{(16)\&(17)}{\leq} \frac{3\sqrt{3}}{2} \left(\frac{K}{\epsilon_k n} \sqrt{\epsilon_k} \max(|\tilde{A}_k|, |\tilde{B}_k|, |\tilde{C}_k|) + \frac{K}{\epsilon_k n} n^{-1/30} \right) \\ &\quad + \frac{\delta}{\epsilon_k n} \max(|\tilde{B}_k|, |\tilde{C}_k|) + \frac{K}{\epsilon_k n} \epsilon_k^{3/4} |\tilde{A}_k| + \frac{K}{\epsilon_k n} n^{-1/30} \\ &\leq \frac{(\delta + 3K\sqrt{\epsilon_k})}{\epsilon_k n} \max(|\tilde{B}_k|, |\tilde{C}_k|) + \frac{4K}{\epsilon_k n} \epsilon_k^{1/2} |\tilde{A}_k| + \frac{4K}{\epsilon_k n} n^{-1/30} \\ &\stackrel{\substack{\text{good region} \\ n \text{ large enough}}}{\leq} K_\epsilon \cdot \left(\frac{2(\delta + 3K\sqrt{\epsilon_k})}{\epsilon_k n} + \frac{4K\epsilon_k^{1/2} |\log \epsilon_k|^{3/4}}{\epsilon_k n} \right). \end{aligned}$$

Up to further diminishing η (which forces $\epsilon_k < \eta$ to be small), we can assume that the right-hand side is smaller than $K_\epsilon/(4\epsilon_k n)$ for n large enough so that we are sure that the conditional expected drift $\mathbb{E} \left[\Delta \tilde{E}_k | \mathcal{F}_k \right]$ is less than $-K_\epsilon/(4\epsilon_k n)$ for $L_E^- < k < L$ and in particular it is negative and pulls back the process towards 0.

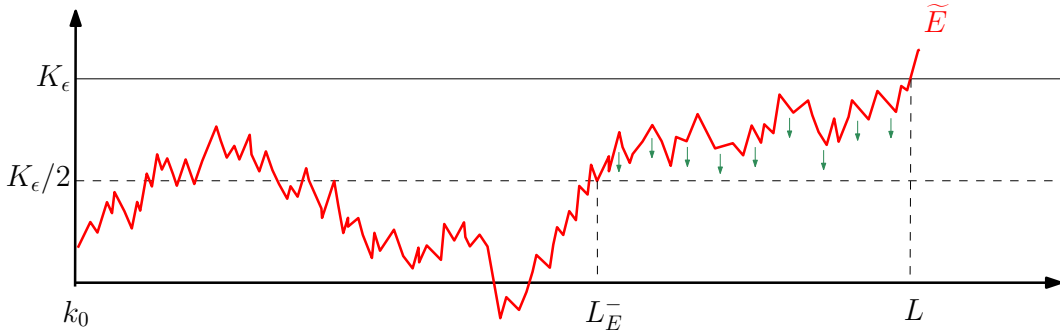


Figure 7: Illustration of the proof. If we exit the good region through the process \tilde{E} , then it has a negative drift (green arrows on the figure) over the time interval (L_E^-, L) and this forces its martingale part to vary too much.

We deduce that on the event $\{k_0 < L_E^- < L < \tilde{\theta} \text{ and } \tilde{E}_L > K_\epsilon\}$ the variation of the martingale M^E over $[L_E^-, L]$ must be larger than $K_\epsilon/2$ (just because the drift plays against the process in this region). Hence,

$$\mathbb{P}(L < \tilde{\theta} \text{ and } \tilde{E}_L > K_\epsilon) \leq \mathbb{P}(L_E^- \leq k_0) + \mathbb{P} \left(\sup_{k_0 \leq k \leq L} |M_{k \wedge L}^E| > \frac{K_\epsilon}{4} \right)$$

We now use our variance estimates (21) and (20) in the good region. In particular,

$$\begin{aligned} \mathbb{E} [(\Delta M_k^E)^2 \mathbb{1}_{k \leq L}] &= \text{Var} (\Delta M_k^E \mathbb{1}_{k \leq L}) = \text{Var} \left(\left(\Delta \tilde{C}_k - \frac{3\sqrt{3}}{2} \Delta \tilde{B}_k \right) \mathbb{1}_{k \leq L} \right) \\ &\leq \text{Cst} \cdot \left(\text{Var}(\Delta \tilde{C}_k \mathbb{1}_{k \leq L}) + \text{Var}(\Delta \tilde{B}_k \mathbb{1}_{k \leq L}) \right) \stackrel{(20) \& (21)}{\leq} \frac{\text{Cst}}{\sqrt{\epsilon_k n}}, \end{aligned}$$

where $\text{Cst} > 0$ as usual does not depend on n nor on ϵ nor on η . Summing those variances over one scale we obtain

$$\sum_{k=k_i}^{k_{i+1}-1} \frac{\text{Cst}}{\sqrt{\epsilon_k n}} \leq \frac{\text{Cst}(k_{i+1} - k_i)}{\sqrt{\epsilon_{k_{i+1}} n}} \leq \text{Cst} \frac{\eta 2^{-i} n}{\sqrt{\eta 2^{-i} n}} = \text{Cst} \sqrt{\eta} (\sqrt{2})^{-i}.$$

We deduce that

$$\mathbb{P} \left(\sup_{k_0 \leq k \leq L} |M_k^E| > \frac{K_\epsilon}{4} \right) \stackrel{\text{Doob}}{\leq} \frac{16 \mathbb{E}[(M_L^E)^2]}{K_\epsilon^2} \leq \frac{\text{Cst}}{K_\epsilon^2} \sum_{i=0}^{\infty} \sum_{k=k_i}^{k_{i+1}-1} \mathbb{E}[(\Delta M_k^E)^2 \mathbb{1}_{k \leq L}] \leq \frac{\text{Cst} \sqrt{\eta}}{K_\epsilon^2}.$$

If K_ϵ is large enough, this bound, as well as $\mathbb{P}(L_E^- \leq k_0)$ (by Lemma 2), are less than ϵ so the probability of the event $\{k_0 < L < \tilde{\theta} \text{ and } \tilde{E}_L > K_\epsilon\}$ is less than 2ϵ . Combined with the symmetric case when $\tilde{E}_L < -K_\epsilon$, this finishes the case of \tilde{E} .

LET'S FINALLY MOVE ON TO THE CONTROL OF \tilde{A} . Again, we decompose \tilde{A} as follows

$$\tilde{A}_k = \tilde{A}_{k_0} + \sum_{\ell=k_0}^{k-1} \mathbb{E} [\Delta \tilde{A}_\ell | \mathcal{F}_\ell] + M_k^A,$$

where $(M_{k \wedge L}^A)_{k \geq k_0}$ is a martingale for the canonical filtration and starts at 0 at time k_0 . Compared to the above cases, we shall look more precisely at the scale of L and introduce

$$J \text{ such that } k_J \leq L < k_{J+1}.$$

In particular, recall that if $k_j \leq k \leq k_{j+1}$ we have $j + |\log_2 \eta| \leq |\log_2 \epsilon_k| \leq (j+1) + |\log_2 \eta|$ so that up to losing a multiplicative factor, we may replace $|\log \epsilon_k|$ by the corresponding scale j in the calculations. As before, let us bound from above the probability that we exit the good region with the process \tilde{A} , that is

$$\mathbb{P}(L < \tilde{\theta} \text{ and } \tilde{A}_L > K_\epsilon \cdot (J+1)^{3/4})$$

and the case $L < \tilde{\theta}$ and $\tilde{A}_L < -K_\epsilon \cdot (J+1)^{3/4}$ is symmetric. As for the case of \tilde{E} , we introduce L_A^- the last time before L where \tilde{A} is smaller than $K_\epsilon(J+1)^{3/4}/2$ and I its corresponding scale (i.e. such that $k_I \leq L_A^- < k_{I+1}$), see Figure 8. As before, we get from Lemma 2 that $L_A^- > k_0$ with high probability when K_ϵ is large. We will now use the fact that the conditional expected drift of \tilde{A} not only goes against \tilde{A} but also that its strength is linear in \tilde{A} .

Specifically, when $L_A^- < k < L$, the process \tilde{A} is larger than $K_\epsilon(J+1)^{3/4}/2$ while the other processes are in absolute value less than K_ϵ so that by (15) the predictable drift is negative and of order $-\tilde{A}_k/(\epsilon_k n)$:

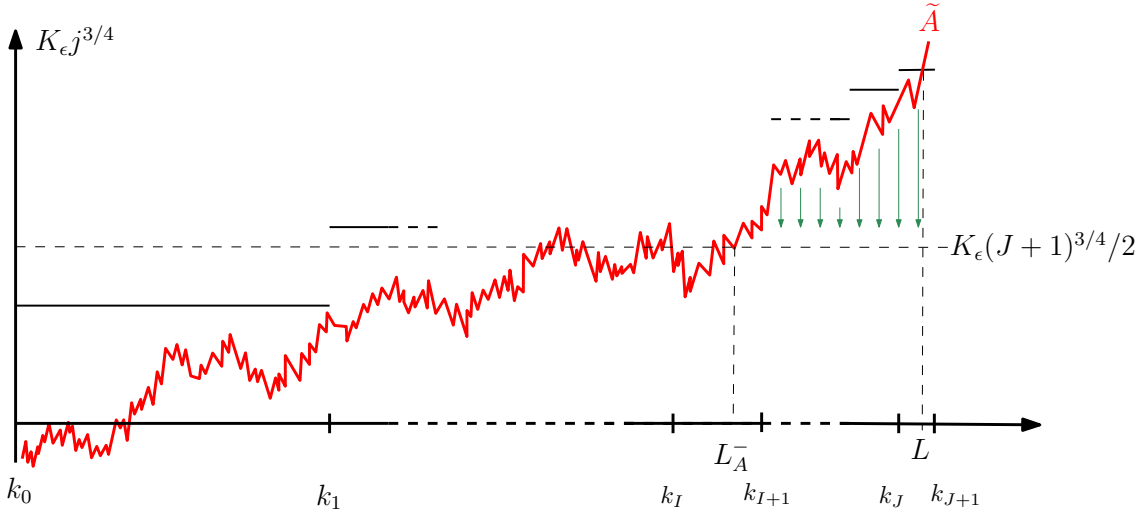


Figure 8: Illustration of the proof. If we exit the good region through the process \tilde{A} , then it has a negative drift (green arrows on the figure) over the time interval (L_A^-, L) whose strength is proportional to $J^{3/4}K_\epsilon$ over a scale. As in the above cases, this forces its martingale part to vary too much.

$$\begin{aligned}
\left| \mathbb{E} \left[\Delta \tilde{A}_k | \mathcal{F}_k \right] - \frac{1}{4\epsilon_k n} \underbrace{\tilde{A}_k}_{\geq K_\epsilon (J+1)^{3/4}/2} \right| &\stackrel{(15)}{\leq} \frac{\delta}{\epsilon_k n} |\tilde{A}_k| + \frac{K\epsilon_k^{1/4}}{\epsilon_k n} \max(|\tilde{B}_k|, |\tilde{C}_k|) + \frac{K}{\epsilon_k n} n^{-1/30} \\
&\stackrel{\text{good region}}{\leq} \frac{1/10}{\epsilon_k n} K_\epsilon (J+1)^{3/4} + \frac{K\epsilon_k^{1/4}}{\epsilon_k n} K_\epsilon + \frac{K}{\epsilon_k n} n^{-1/30} \\
&\leq \frac{1}{9\epsilon_k n} K_\epsilon (J+1)^{3/4}, \tag{38}
\end{aligned}$$

for n large enough up to diminishing η if necessary. In particular, $\mathbb{E} [\Delta \tilde{A}_k | \mathcal{F}_k]$ is less than $-\frac{K_\epsilon}{100} \frac{(J+1)^{3/4}}{\epsilon_k n}$ and summing the conditional expected drift over all $k \in (L_A^-, L)$ yields total drift smaller than

$$\sum_{k=L_A^-+1}^{L-1} \mathbb{E}[\Delta A_k | \mathcal{F}_k] \leq -c \cdot (J - I - 1) K_\epsilon (J+1)^{3/4},$$

for some constant $c > 0$. Let us first concentrate on the case where $I + 1 < J$ so that $J - I - 1 > 0$. In particular, the variation of the martingale M^A between $L_A^- + 1$ and L must compensate this drift and must be larger than $-c(J - I - 1)K_\epsilon(J + 1)^{3/4}$. Thus, we have

$$\begin{aligned}
&\mathbb{P} \left(k_0 < L_A^- < L < \tilde{\theta} \text{ and } I + 1 < J \text{ and } \tilde{A}_L / (J + 1)^{3/4} > K_\epsilon \right) \\
&\leq \sum_{j=2}^{\log_2(n)} \sum_{i=0}^{j-2} \mathbb{P} \left(\sup_{k_i \leq \ell < k_{j+1} \wedge L} M_\ell^A - \inf_{k_i \leq \ell < k_{j+1} \wedge L} M_\ell^A > c(j - i - 1) K_\epsilon (j + 1)^{3/4} \right) \\
&\stackrel{D_{\text{dob}}}{\leq} \text{Cst} \sum_{j=2}^{\log_2(n)} \sum_{i=0}^{j-2} \frac{\mathbb{E}[(M_{k_{j+1} \wedge L}^A - M_{k_i}^A)^2]}{K_\epsilon^2 (j + 1)^{3/2} (j - i - 1)^2}.
\end{aligned}$$

Thanks to our variance estimates (19) we have $\mathbb{E}[(\Delta M_k^A)^2 \mathbf{1}_{k < L}] \leq \frac{\text{Cst}}{\epsilon_k n}$ so that after summing over scales we obtain $\mathbb{E}[(M_{k_{j+1} \wedge L}^A - M_{k_i}^A)^2] \leq \text{Cst} \cdot (j + 1 - i)$. Plugging this back into the above estimate

we deduce

$$\begin{aligned} & \mathbb{P} \left(k_0 < L_A^- < L < \tilde{\theta} \text{ and } I + 1 < J \text{ and } \tilde{A}_L / (J + 1)^{3/4} > K_\epsilon \right) \\ & \leq \frac{\text{Cst}}{K_\epsilon^2} \sum_{j=2}^{\log_2(n)} \sum_{i=0}^{j-2} \frac{(j+1-i)}{(j+1)^{3/2}(j-i-1)^2} \leq \frac{\text{Cst}}{K_\epsilon^2}, \end{aligned}$$

so that this probability can be made arbitrarily small by making K_ϵ large. The case $\tilde{A}_L / (J + 1)^{3/4} < -K_\epsilon$ is treated similarly. As for the case $|I - J| \leq 1$, since $k_0 < L_A^-$ w.h.p. (by Lemma 2), we use that in this case the martingale M^A must have a variation of at least $K_\epsilon(j+1)^{3/4}/2$ over (k_{j-1}, k_{j+1}) (we do not use the strength of the drift, but just the fact it plays against us over (L_A^-, L) as for the case of \tilde{E}). By Doob maximal inequality and the above estimate, this probability is bounded from above by $\text{Cst}/((j+1)^{3/2}K_\epsilon^2)$, whose sum over $0 \leq j \leq \log_2(n)$ is $\leq \frac{\text{Cst}}{K_\epsilon^2}$. We conclude similarly. \square

4.3 This is the end

Using Proposition 6 and (8), we can conclude as in Lemma 1 that the process X_k stays positive at least as long as

$$n\varepsilon_k^2 \gg \sqrt{n}(\varepsilon_k)^{3/4} |\log \varepsilon_k|^{3/4}, \quad \text{i.e. as long as } t_{\text{ext}}n - k \gg n^{3/5}(\log n)^{3/5}.$$

Through a more refined control on \tilde{A} , we shall first prove that we can remove the $\log^{3/5} n$ and prove that $nt_{\text{ext}} - \theta = O_{\mathbb{P}}(n^{3/5})$, see Proposition 7. The convergence in law of $n^{-3/5}(nt_{\text{ext}} - \theta)$ will be deduced by doing a SDE approximation for the process \tilde{A}_k when $\varepsilon_k \approx n^{-2/5}$ in Proposition 8.

Since we now take a close look at times $k = t_{\text{ext}}n - O(n^{3/5})$, let us introduce a new piece of notation: for $k \geq 0$, we write

$$\mathbf{t}_k := n^{-3/5}(t_{\text{ext}}n - k), \quad \text{so that } k = t_{\text{ext}}n - \mathbf{t}_k n^{3/5} \quad \text{i.e. } \varepsilon_k = \mathbf{t}_k n^{-2/5}.$$

With this notation at hands, we can state a refined control on \tilde{A} .

Proposition 7 (Control on \tilde{A} in the critical region). *For all $\epsilon > 0$ there exists K_ϵ such that with probability at least $1 - \epsilon$, for all $k \leq t_{\text{ext}}n$ such that $\mathbf{t}_k \geq K_\epsilon$, we have*

$$|\tilde{A}_k| < K_\epsilon \mathbf{t}_k^{1/8}.$$

In particular, performing the same argument as in the beginning of this subsection, we deduce that X_k stays positive until time $t_{\text{ext}}n - K_\epsilon n^{3/5}$, that is $\theta > t_{\text{ext}}n - K_\epsilon n^{3/5}$ with probability at least $1 - \epsilon$.

Proof. The proof is similar to the control of \tilde{A} in Proposition 6. With the notation of the proof of Proposition 6, let us introduce

$$T = L \wedge \min\{k \geq k_0 : |\tilde{A}_k| > K_\epsilon \cdot \mathbf{t}_k^{\frac{1}{8}}\}$$

and $J \in \{0, 1, 2, \dots\}$ the corresponding scale, i.e. such that $2^J \geq \mathbf{t}_T > 2^{(J-1)}$. As for the previous control of \tilde{A} , we will replace \mathbf{t}_T by 2^J in the calculation to make the reading easier. Note in particular that $k \mapsto \mathbf{t}_k$ is decreasing.

We bound the probability $\mathbb{P}(T = k < L \text{ and } \tilde{A}_k > K_\epsilon \cdot 2^{\frac{J}{8}})$, the case $\{T = k < L \text{ and } \tilde{A}_k < -K_\epsilon \cdot 2^{\frac{J}{8}}\}$ being similar. For this, let $\alpha > 0$ be a small constant (to be precised later), and let

$$T^- = \sup \left\{ k_0 \leq k \leq T : \tilde{A}_k \leq \alpha(I - J + 1)K_\epsilon 2^{J/8} \text{ with } 2^{I-1} < \mathbf{t}_k \leq 2^I \right\}$$

and $I \geq J$ its corresponding scale (notice the slight difference here with the proof of Proposition 6 because I enters in the definition of the barrier). As before, Lemma 2 will entail that $T^- > k_0$ with high probability as $n \rightarrow \infty$ and when $k_0 \leq T^- \leq k \leq T$ and $2^{i-1} < \mathbf{t}_k \leq 2^i$, we have $\tilde{A}_k \geq \alpha(i - J + 1)K_\epsilon 2^{J/8}$. By the same calculation as in (38) we have

$$\mathbb{E} \left[\Delta \tilde{A}_k | \mathcal{F}_k \right] \leq -\frac{\alpha(i - J + 1)K_\epsilon 2^{J/8}}{8 \epsilon_k n}.$$

Summing those expected conditional drifts over all $T^- + 1 \leq k < T$ yields a total drift smaller than

$$\begin{aligned} \sum_{k=T^-+1}^{T-1} \mathbb{E}[\Delta \tilde{A}_k | \mathcal{F}_k] &\leq \sum_{\text{scales}} \sum_{i=J+1}^{I-1} \sum_{k \geq 0} \mathbb{1}_{2^i \geq \mathbf{t}_k > 2^{i-1}} \mathbb{E}[\Delta \tilde{A}_k | \mathcal{F}_k] \\ &\leq -\frac{\alpha}{8} \sum_{i=J+1}^{I-1} (i - J + 1) K_\epsilon 2^{J/8} \sum_{k \geq 0} \mathbb{1}_{2^i \geq \mathbf{t}_k > 2^{i-1}} \frac{1}{\epsilon_k n} \\ &\leq -\frac{\alpha}{16} \sum_{i=J+1}^{I-1} (i - J + 1) K_\epsilon 2^{J/8} \\ &\leq -\frac{\alpha}{16} (I - J - 1)^2 K_\epsilon 2^{J/8}. \end{aligned}$$

Let us first focus on the case $I - J \geq 2$: as soon as $T^- > k_0$ the variation of the martingale M^A between T^- and T must compensate this drift plus the difference of the starting and ending values, and so must be larger than

$$K_\epsilon 2^{J/8} \left(\frac{\alpha}{16} (I - J - 1)^2 - \alpha(I - J + 1) + 2^{-1/8} \right).$$

If α has been chosen small enough (e.g. $\alpha = \frac{1}{100}$), as soon as $I - J \geq 2$, this is larger $\frac{1}{32} K_\epsilon 2^{J/8} (I - J)^2$. As in the proof of Proposition 6, the sum of the variances of the increments of M^A between scales i and j is bounded above by $\text{Cst}(i - j)$ and so the probability that M^A varies by more than $\frac{1}{32}(i - j)^2 K_\epsilon 2^{j/8}$ over this time interval is bounded above using Doob's inequality by

$$\text{Cst} \frac{i - j}{((i - j)^2 K_\epsilon 2^{j/8})^2}.$$

Summing these probabilities over all scales $j_0 \leq j \leq i$, we deduce that

$$\begin{aligned} &\mathbb{P} \left(k_0 < T^- < T < L \text{ and } I - 1 > J \geq j_0 \text{ and } \tilde{A}_T > K_\epsilon 2^{J/8} \right) \\ &\leq \frac{\text{Cst}}{K_\epsilon^2} \sum_{i \geq j+2 \geq j_0+2} \frac{i - j}{(i - j)^4 2^{j/4}} \leq \frac{\text{Cst} \cdot 2^{-j_0/4}}{K_\epsilon^2}, \end{aligned}$$

and this can be made arbitrarily small by taking j_0 large enough. Finally, we treat the case $0 \leq I - J \leq 1$ similarly, by noting that in this case, if $k_0 < T^-$ (which has high probability by Lemma 2), the variation of \tilde{A} between times T^- and T is at least $(2^{-1/8} - 2\alpha)K_\epsilon 2^{J/8}$. Since the drift is negative, the martingale M^A must have a variation of order $K_\epsilon 2^{J/8}$ (provided $\alpha < \frac{1}{4}$) over the scale J , and the conclusion is the same. \square

In the rest of this subsection we stress back the dependence in n and use $\theta^n \equiv \theta$ for the stopping time of the exploration and study the convergence of

$$\mathbf{t}_{\theta^n} \in \mathbb{R} \quad \text{such that} \quad \theta^n = t_{\text{ext}} n - \mathbf{t}_{\theta^n} n^{3/5}.$$

Proposition 8. *We have the following convergence in distribution as n goes to infinity*

$$\mathbf{t}_{\theta^n} \xrightarrow[n \rightarrow \infty]{(d)} 3^{-3/5} \cdot 2^{4/5} \cdot \vartheta^{-2},$$

where $\vartheta = \inf\{t \geq 0 : W_t = -t^{-2}\}$ with W a standard linear Brownian motion started from 0 at 0.

Proof. Fix $\epsilon > 0$ and let $K_\epsilon > 0$ so that on an event \mathcal{E}_n of probability at least $1 - 3\epsilon$, the conclusions of Lemma 2, Proposition 7 and Proposition 6 hold. Fix $K_\epsilon^{-1} > \xi > 0$ small enough so that $K_\epsilon \xi^{1/8} \leq \epsilon$. We shall first focus on the times k satisfying $\xi \leq \mathbf{t}_k \leq \xi^{-1}$ and consider the renormalized process

$$\tilde{F}_k = \frac{\tilde{A}_k}{\mathbf{t}_k^{1/4}}, \quad 0 \leq k \leq \theta^n.$$

Let us compute its conditional expected drift and variance: for $k < \tilde{\theta}^n$ with $\xi \leq \mathbf{t}_k \leq \xi^{-1}$, on the event \mathcal{E}_n the assumptions of Proposition 4 hold, so that using $\varepsilon_k = n^{-2/5} \mathbf{t}_k$ we have

$$\mathbb{E}[\Delta \tilde{F}_k | \mathcal{F}_k, \mathcal{E}_n] \leq \frac{\delta}{\mathbf{t}_k n^{3/5}} |\tilde{A}_k| + \frac{K}{\mathbf{t}_k n^{3/5}} n^{-1/30} = \frac{\delta}{\mathbf{t}_k^{3/4} n^{3/5}} |\tilde{F}_k| + \frac{K}{\mathbf{t}_k n^{3/5}} n^{-1/30} \quad (39)$$

$$\left| \text{Var} \left(\Delta \tilde{F}_k | \mathcal{F}_k, \mathcal{E}_n \right) - \frac{2\sqrt{3}}{\mathbf{t}_k^{3/2} n^{3/5}} \right| = \left| \frac{1}{\mathbf{t}_k^{1/2}} \text{Var} \left(\Delta \tilde{A}_k | \mathcal{F}_k \right) - \frac{2\sqrt{3}}{\mathbf{t}_k^{3/2} n^{3/5}} \right| \leq \frac{\delta}{\mathbf{t}_k^{3/2} n^{3/5}}. \quad (40)$$

We now make δ vary with n and take $\delta \equiv \delta_n \xrightarrow[n \rightarrow \infty]{} 0$ in the above displays. Indeed, using the notation of Propositions 4 and 5 we can do so as soon as $\eta(\delta_n) > 1/\xi \cdot n^{-2/5}$. To avoid stopping times issues, we possibly extend \tilde{F} after time θ^n (in the case $\mathbf{t}_{\theta^n} \leq \xi$) by a process \hat{F} whose increments are $\pm (\frac{2\sqrt{3}}{\mathbf{t}_k^{3/2} n^{3/5}})^{1/2}$ with probability 1/2 (in particular independent, centered, with variance $\frac{2\sqrt{3}}{\mathbf{t}_k^{3/2} n^{3/5}}$ and whose L^∞ -norm tends to 0 uniformly as $n \rightarrow \infty$), so that our estimates (39) and (40) remain true for all $\{k : \xi \leq \mathbf{t}_k \leq \xi^{-1}\}$. Let us recapitulate what we have: with probability at least $1 - 3\epsilon$ for all $\{k : \xi \leq \mathbf{t}_k \leq \xi^{-1}\}$:

$$\left\{ \begin{array}{l} |\hat{F}_{n t_{\text{ext}} - \xi^{-1} n^{3/5}}| < \epsilon, \quad (\text{by Prop. 7 and the assumption } K_\epsilon \xi^{1/8} \leq \epsilon), \\ \mathbb{E}[\Delta \hat{F}_k | \mathcal{F}_k] = o(n^{-3/5}) \cdot |\hat{F}_k| + o(n^{-3/5}), \\ \text{Var} \left(\Delta \hat{F}_k | \mathcal{F}_k \right) = \frac{2\sqrt{3}}{\mathbf{t}_k^{3/2} n^{3/5}} + o(n^{-3/5}), \\ \|\Delta \hat{F}_k\|_\infty = o(1), \end{array} \right.$$

where the $o(1)$ function is uniform in $\{k : \xi \leq \mathbf{t}_k \leq \xi^{-1}\}$. By standard results in diffusion approximation, see e.g. [14], this implies the following weak convergence for the $\|\cdot\|_\infty$ -norm:

$$\left(\hat{F}_{t_{\text{ext}} n - t n^{3/5}} - \hat{F}_{t_{\text{ext}} n - \xi^{-1} n^{3/5}} \right)_{\xi \leq t \leq \xi^{-1}} \xrightarrow[n \rightarrow \infty]{} (\mathcal{H}_t)_{\xi \leq t \leq \xi^{-1}},$$

where the process \mathcal{H} satisfies the stochastic differential equation (in reverse time) $d\mathcal{H}_{-t} = \frac{\sqrt{2\sqrt{3}}}{t^{3/4}} dB_{-t}$ with initial condition $\mathcal{H}_{\xi^{-1}} = 0$. By Dubbins-Schwarz theorem, the solution of this SDE can be written as

$$2 \cdot 3^{1/4} \left(W_{\frac{1}{\sqrt{t}}} - W_{\frac{1}{\sqrt{\xi^{-1}}}} \right)_{\xi \leq t \leq \xi^{-1}}$$

where W is a standard linear Brownian motion with $W_0 = 0$. Letting $\epsilon \rightarrow 0$ and $\xi \rightarrow 0$, we deduce the following convergence weak convergence over all compact subsets of $(0, \infty)$:

$$\left(\widehat{F}_{t_{\text{ext}}n - tn^{3/5}} \right)_{0 < t < \infty} \xrightarrow[n \rightarrow \infty]{} \left(W_{\frac{1}{\sqrt{t}}} \right)_{0 < t < \infty}. \quad (41)$$

To see that the above convergence implies the convergence of stopping times recall that

$$\begin{aligned} \mathbf{t}_{\theta^n} := \sup\{\mathbf{t}_k \geq 0, X_k = 0\} &= \sup\{\mathbf{t}_k \geq 0, \widetilde{F}_k = -n^{4/5} \mathcal{X}(k/n)/\mathbf{t}_k\} \\ &= \sup\{\mathbf{t}_k \geq 0, \widehat{F}_k \leq -n^{4/5} \mathcal{X}(k/n)/\mathbf{t}_k\}. \end{aligned}$$

In particular, the time \mathbf{t}_{θ^n} can be seen as the first time when started from $+\infty$ that the process \widehat{F} crosses the barrier \mathcal{C}^n defined by

$$\mathcal{C}^n(\mathbf{t}_k) = -n^{4/5} \mathcal{X}(k/n)/\mathbf{t}_k.$$

Recalling (8), we have $-n^{4/5} \mathcal{X}(k/n)/\mathbf{t}_k \sim -3\mathbf{t}_k$, so that the barrier \mathcal{C}^n converges towards the graph \mathcal{C} of the function $t \mapsto -3t$. Since the crossing of \mathcal{C} by $(W_{1/\sqrt{t}} : 0 < t < \infty)$ when started from $+\infty$ happens at an almost surely positive time τ and since W immediately takes values strictly above and below \mathcal{C} after hitting it, it follows that

$$\mathbf{t}_{\theta^n} \xrightarrow[n \rightarrow \infty]{(d)} \tau = \sup\{t \geq 0 : 2 \cdot 3^{1/4} \cdot W_{\frac{1}{\sqrt{t}}} = -3t\}.$$

By scaling we have the equality in distribution

$$\begin{aligned} \tau &\stackrel{(d)}{=} \sup\{t \geq 0 : 2 \cdot 3^{1/4} \cdot W_{\frac{1}{\sqrt{t}}} = -3t\} \\ &=_{u=1/\sqrt{t}} \left(\inf\{u \geq 0 : W_u = \frac{3^{3/4}}{2} u^{-2}\} \right)^{-2} \\ &\stackrel{(d)}{=}_{\alpha > 0} \left(\inf\{u \geq 0 : \frac{1}{\sqrt{\alpha}} \cdot W_{\alpha u} = \frac{3^{3/4}}{2} u^{-2}\} \right)^{-2} \\ &=_{\alpha u = v} \left(\frac{1}{\alpha} \inf\{v \geq 0 : W_v = \sqrt{\alpha} \alpha^2 \cdot \frac{3^{3/4}}{2} v^{-2}\} \right)^{-2} \\ &=_{\alpha^{5/2} \cdot \frac{3^{3/4}}{2} = 1} \left(\frac{3^{3/4}}{2} \right)^{-4/5} \left(\inf\{v \geq 0 : W_v = v^{-2}\} \right)^{-2}. \end{aligned}$$

The statement follows. \square

4.4 Proof of Theorem 2: Size and composition of the KS-Core

We have now all the ingredients to prove our main Theorem 2. First by Proposition 8, the renormalized ending time \mathbf{t}_{θ^n} converges in distribution to $2^{4/5} 3^{-3/5} \vartheta^{-2}$ where ϑ is the hitting time of the curve $t \mapsto -t^{-2}$ by a Brownian motion. At this time, by Proposition 6 and (8) we have

$$Y_{\theta^n} = \underbrace{B_{\theta^n}}_{\stackrel{\leq}{\text{Prop. 6}} \text{ Cst} \sqrt{n} \log(n)^{3/4}} + \underbrace{n \mathcal{Y} \left(\frac{\theta^n}{n} \right)}_{\stackrel{\sim}{(8)} 4 \mathbf{t}_{\theta^n} n^{3/5}} \approx 4 \mathbf{t}_{\theta^n} n^{3/5},$$

$$Z_{\theta^n} = \underbrace{C_{\theta^n}}_{\substack{\leq \text{Cst } n^{3/10} \log(n)^{3/4} \\ \text{Prop. 6}}} + \underbrace{n \mathcal{L}\left(\frac{\theta^n}{n}\right)}_{\substack{\sim 4\sqrt{3}t_{\theta^n} n^{2/5} \\ (8)}} \approx 4\sqrt{3}t_{\theta^n}^{3/2} n^{2/5}.$$

Moreover using Proposition 1, the KS-Core is just obtained by pairing the remaining half-edges uniformly at random. Our theorem follows. Ouff.

5 Comments

We conclude this paper with a few perspectives that our work opens.

Near critical heuristics. We believe that our techniques can be used to tackle the near-critical window for the Karp-Sipser core. In particular, this window should be obtained by starting from

$$d_{1,c}^n = n\left(1 - \frac{\sqrt{3}}{2}\right) + O(n^{3/5}), \quad 2d_{2,c}^n = O(n^{3/5}), \quad \text{and} \quad 3d_{3,c}^n = n\frac{\sqrt{3}}{2} + O(n^{3/5}),$$

whereas we studied only the critical case (3). All these shifts in the starting configuration should result in a shift of order $O(n^{3/5})$ of the absorption time. In a similar vein, one could study the “Phase 2” of the Karp-Sipser algorithm [1] which, in the supercritical case, consists in removing a uniform vertex when there are no leaves left. The analysis of this phase should be intimately connected to the above near-critical dynamics.

Universality. Obviously, we conjecture that the geometry of the critical core and the scaling limits results are independent of the fine details of the model of random graph we started with. In particular, it should hold for the Erdős-Rényi case or for configuration models with small enough degrees. However, proving a general result seems challenging because we heavily rely on the exact form of the fluid-limit of our exploration processes (such results are available for the Erdős-Rényi case, see [1]).

Stopped Markov chain. More generally, we believe that the techniques developed in this paper could be used to understand precisely the exit times of Markov chains from domains. To fix ideas, let $(\mathbb{X}^n : k \geq 0)$ be a \mathbb{Z}^d -valued Markov chain whose expected conditional drift is well-approximated by $\phi(\mathbb{X}^n/n)$ for some function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^d$. The differential equation method shows that under some mild assumptions $(n^{-1}\mathbb{X}_{\lfloor tn \rfloor}^n : t \geq 0)$ converges towards a solution \mathcal{X} to $\mathcal{X}'(t) = \phi(\mathcal{X}(t))$. If Ω is a bounded domain and Ω^n its discrete approximation, it is reasonable to believe that the exist time θ^n of Ω^n by \mathbb{X}^n should converge after normalization towards the exit time t_{ext} of Ω by \mathcal{X} . However, the fine fluctuations of θ^n around nt_{ext} should depend on fine properties of ϕ (and its derivatives) near the exit point. We plan on addressing those general questions in future works.

Comparison with the k -core phase transition. Finally, it is interesting to compare our results with the appearance of the k -core in random graphs as studied in [16, 9], where the phase transition is discontinuous.

Recall that the k -core of a graph \mathfrak{g} is the maximal subgraph of $\mathfrak{g}' \subset \mathfrak{g}$ so that the induced degree inside \mathfrak{g}' of each of its vertices is at least k . The emergence of a giant k -core has been studied for

the Erdős–Rényi random graph and the configuration model, see [16, 9]. A difference with the Karp–Sipser core is that the phase transition is discontinuous: when the k -core exists asymptotically, its proportion is bounded away from 0. This can be explained heuristically as follows.

Suppose for the discussion that $k = 3$ and that we are interested in the size of the 3-core in a configuration model on vertices of degrees 1, 2, 3 and 4. As in the case of the Karp–Sipser algorithm, one can reveal the 3-core by iteratively taking a leg attached to a vertex of degree ≤ 2 , remove it, and destroy the vertex it is attached to as well as the connection it makes in the graph (hence diminishing the unmatched degree of the vertices in question). As in this paper, if one starts with some proportions p_1, p_2, p_3, p_4 of legs attached to vertices of degree one, two, three and four, we can write the differential equation governing the fluid limit of this process, see [9]. The main difference with the Karp–Sipser core is that in this case, the number of legs attached to leaves (to be precise to vertices of degree 1 or 2) is not necessarily decreasing. Actually, in the critical case, the fluid limit of the proportion of vertices of degrees 1, 2 follows a curve which is tangent to the boundary of the domain at some point before diving back into the bulk of the simplexe and dying at the right corner, see Figure 9 (and compare with Figure 5). This explains the first-order phase transition in this case: a slight perturbation of the initial conditions may push the curve to exit the domain at a very different location.

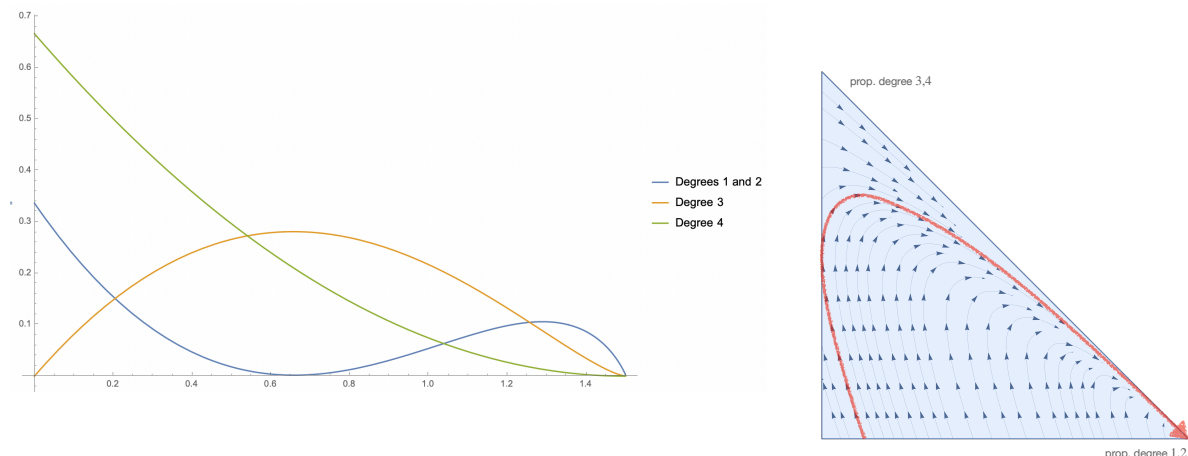


Figure 9: Illustration of the fluid limit of the renormalized number of legs attached to vertices of degree 4, 3 and ≤ 2 in the k -core algorithm at the critical point. In particular, a slight perturbation of the initial conditions may cause a drastic change of the absorption time of the system and this explains why the k -core percolation exhibits a first-order phase transition.

References

- [1] J. ARONSON, A. FRIEZE, AND B. G. PITTEL, *Maximum matchings in sparse random graphs: Karp–sipser revisited*, Random Structures & Algorithms, 12 (1998), pp. 111–177.

- [2] M. BAUER AND O. GOLINELLI, *Core percolation in random graphs: a critical phenomena analysis*, The European Physical Journal B-Condensed Matter and Complex Systems, 24 (2001), pp. 339–352.
- [3] M. BAUER AND O. GOLINELLI, *Random incidence matrices: moments of the spectral density*, Journal of Statistical Physics, 103 (2001), pp. 301–337.
- [4] T. BOHMAN AND A. FRIEZE, *Karp-sipser on random graphs with a fixed degree sequence*, Combinatorics, Probability & Computing, 20 (2011), pp. 721–741.
- [5] B. BOLLOBÁS, *A probabilistic proof of an asymptotic formula for the number of labelled regular graphs*, European J. Combin., 1 (1980), pp. 311–316.
- [6] S. COSTE AND J. SALEZ, *Emergence of extended states at zero in the spectrum of sparse random graphs*, The Annals of Probability, 49 (2021), pp. 2012–2030.
- [7] S. N. ETHIER AND T. G. KURTZ, *Markov processes: characterization and convergence*, John Wiley & Sons, 2009.
- [8] C. GOLDSCHMIDT AND E. KREAČIĆ, *The spread of fire on a random multigraph*, Advances in Applied Probability, 51 (2019), pp. 1–40.
- [9] S. JANSON AND M. J. LUCZAK, *A simple solution to the k -core problem*, Random Structures & Algorithms, 30 (2007), pp. 50–62.
- [10] M. JONCKHEERE AND M. SÁENZ, *Asymptotic optimality of degree-greedy discovering of independent sets in configuration model graphs*, Stochastic Processes and their Applications, 131 (2021), pp. 122–150.
- [11] R. M. KARP AND M. SIPSER, *Maximum matching in sparse random graphs*, in 22nd Annual Symposium on Foundations of Computer Science (sfcs 1981), IEEE, 1981, pp. 364–375.
- [12] E. KREACIC, *Some problems related to the Karp-Sipser algorithm on random graphs*, PhD thesis, University of Oxford, 2017.
- [13] T. G. KURTZ, *Solutions of ordinary differential equations as limits of pure jump markov processes*, Journal of applied Probability, 7 (1970), pp. 49–58.
- [14] H. J. KUSHNER, *On the weak convergence of interpolated markov chains to a diffusion*, The annals of Probability, (1974), pp. 40–50.
- [15] M. MOLLOY AND B. REED, *A critical point for random graphs with a given degree sequence*, Random structures & algorithms, 6 (1995), pp. 161–180.
- [16] B. PITTEL, J. SPENCER, AND N. WORMALD, *Sudden emergence of a giant k -core in a random graph*, Journal of Combinatorial Theory, Series B, 67 (1996), pp. 111–151.
- [17] J. SALEZ, *Some implications of local weak convergence for sparse random graphs*, PhD thesis, Université Pierre et Marie Curie-Paris VI; Ecole Normale Supérieure de Paris . . . , 2011.
- [18] R. VAN DER HOFSTAD, *Random graphs and complex networks. vol. ii*, available at <http://www.win.tue.nl/~rhofstad/>, (preliminary version).
- [19] N. C. WORMALD, *Differential equations for random processes and random graphs*, The annals of applied probability, (1995), pp. 1217–1235.