

Feedforward Mixing is as Sharp as it is Slow in Reverse

Anonymous Authors¹

Abstract

Deep learning architectures across diverse domains—including language modeling, audio-visual processing, and temporal graph modeling—fundamentally rely on feedforward computational graphs. Consequently, optimising these graph topologies is a major research focus. Recent theoretical work evaluates the efficacy of these graphs using two key metrics: forward mixing time $T_{fwd}(G)$ (the speed of information propagation) and normalised minimax fidelity $F(G)$ (the sharpness of information representation). In this work, we prove that, surprisingly, minimax fidelity is tightly upper-bounded by the *backward* mixing time of computational graphs: $F(G) \leq \frac{8}{3} \cdot T_{bwd}(G)$ for all feedforward graphs with a unique source and sink. Practically, this establishes a direct tradeoff, $F(G) \leq \frac{8}{3} \cdot T_{fwd}(G)$, for symmetric graphs where $T_{fwd}(G) = T_{bwd}(G)$. In the uniform setting, because most common architectures—including fully-connected and sliding-window graphs—are symmetric, this exposes an unavoidable architectural limit: optimising for rapid information mixing necessitates a degradation in fidelity, and vice versa. More generally, i.e. if $T_{bwd}(G) = f(T_{fwd}(G))$ and consequently $F(G) \leq f(T_{fwd}(G))$ for some function f , we show that f is well-behaved and at most an $O(n)$ factor of $T_{fwd}(G)$ for both uniform and non-uniform graphs, where n is the number of nodes in G . Our results, ultimately, establish backward mixing time as a key metric for evaluating information propagation in deep learning models.

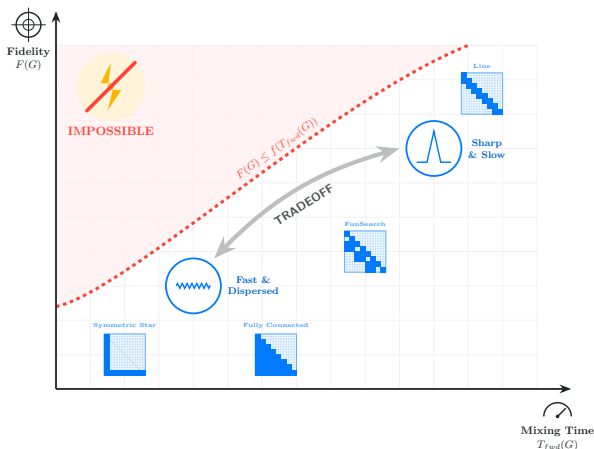


Figure 1. More generally, the normalised minimax fidelity $F(G)$ is upper-bounded by a function of the mixing time $T_{fwd}(G)$. This implies that, across the design space dictated by minimax fidelity and mixing time, feedforward computational graphs lie on or below the red boundary.

1. Introduction

Many problems in deep learning necessitate the processing of sequential inputs (Vitvitskyi et al., 2025). Language modelling (Sutskever et al., 2014), audio and video processing (Wiedemer et al., 2025; van den Oord et al., 2016), time-series forecasting (Kim et al., 2025), and temporal graph modelling (Huang et al., 2023) all employ causal processing of their inputs, and at their core, the forward passes of these sequence models can be abstracted as feedforward computational graphs. Significant research effort has been dedicated to optimising computational graph topologies to facilitate information flow, albeit most of the effort has happened in the *undirected* graph space. Notable examples include rewiring strategies in graph neural networks (GNNs) to avoid over-squashing and under-reaching (Attali et al., 2024; Giovanni et al., 2023; Topping et al., 2022; Deac et al., 2022; Wilson et al., 2024), and the development of dy-

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

055 namic, data-dependent attentional masks (Xia et al., 2022)
 056 for vision Transformers. In the feedforward case, several
 057 works have designed static attentional masks for decoder-
 058 only Transformers (Beltagy et al., 2020; Chen et al., 2025;
 059 Child et al., 2019).

060 Given the ubiquity of computational feedforward graphs,
 061 understanding the theoretical limits of their design is
 062 paramount. To formally analyse this design space, (Vitvit-
 063 skyi et al., 2025) introduced two metrics to quantify the
 064 efficacy of information propagation: forward *mixing time*
 065 (how rapidly information traverses the computational graph;
 066 see Espuny Dfáz et al. (2024)) and normalised minimax
 067 *fidelity* (how distinctly information from any source node is
 068 preserved at the sink). Intuitively, an ideal computational
 069 graph would simultaneously achieve **low mixing time** (fast
 070 global propagation) and **high minimax fidelity** (sharp, dis-
 071 tinguishable representations).
 072

073 In this paper, we prove that normalised minimax fidelity
 074 is tightly upper-bounded by the **backward** mixing time of
 075 computational graphs: $F(G) \leq \frac{8}{3} \cdot T_{bwd}(G)$ for all feedfor-
 076 ward graphs with a unique source and a unique sink. Practi-
 077 cally, this establishes a direct tradeoff, $F(G) \leq \frac{8}{3} \cdot T_{fwd}(G)$,
 078 specifically for graphs that are both **symmetric** (equivalent
 079 up to isomorphism when reversed; Section 5.2) and **uni-**
 080 **form** (where edge weights are assigned uniformly based
 081 on in-degree and out-degree). For symmetric uniform
 082 graphs, which include fully-connected and sliding-window
 083 graphs, this exposes an unavoidable architectural limit: op-
 084 timising for rapid information mixing necessitates a degra-
 085 dation in fidelity, and vice versa. In the more general
 086 case, i.e. if $T_{bwd}(G) = f(T_{fwd}(G))$ and consequently
 087 $F(G) \leq f(T_{fwd}(G))$ for some function f , we show that
 088 $T_{bwd}(G)$ and $T_{fwd}(G)$ diverge by at most an $O(n)$ factor
 089 for uniform graphs and $O(n/c)$ for non-uniform graphs
 090 whose weights are bounded away from zero by c . With
 091 these findings, our results establish backward mixing time
 092 as a key metric for evaluating information propagation. We
 093 summarise our core contributions below:
 094

- 095 • **A Duality Between Fidelity and Backward Mixing.**
 096 We first establish an equivalence between the forward
 097 diffusion of representations and a random walk on the
 098 reversed computational graph.
- 099 • **Fidelity is Tightly Bounded by the Backward Mix-**
 100 **ing Time.** Using this duality, we prove our central re-
 101 sult: minimax fidelity is tightly upper-bounded by the
 102 graph’s worst-case backward mixing time for graphs
 103 with a unique source and a unique sink. Furthermore,
 104 we identify specific topologies (conjoined star graphs)
 105 to demonstrate that this bound is tight.
- 106 • **Limits on Forward Mixing and Fidelity.** We connect
 107 our primary bound back to the forward mixing time
 108
 109

$T_{fwd}(G)$. We make the observation that for symmetric
 graphs—encompassing a wide range of practical
 architectures—forward and backward mixing times
 are identical, exposing a tight $F(G) \leq \frac{8}{3} \cdot T_{fwd}(G)$
 tradeoff. More generally, $F(G) \leq f(T_{fwd}(G))$ holds,
 where $f(x) = O(n) \cdot x$ and $f(x) = O(n/c) \cdot x$ for
 uniform and non-uniform graphs respectively, where n
 is the size of the graph G .

- **Empirical Demonstration.** We plot the worst-case
 forward mixing time and minimax fidelity across a
 wide range of standard attentional masks in the uni-
 form graph setting, demonstrating that our bounds hold
 empirically.

2. Related Work

The theoretical framing of information flow in deep learn-
 ing architectures has gained traction as a means to explain
 empirical bottlenecks. Our work mainly builds upon (Vitvit-
 skyi et al., 2025), who introduced mixing time and mini-
 max fidelity as rigorous metrics to quantify the propagation
 characteristics of feedforward computational graphs. Build-
 ing on this framework, (de Ocariz Borde, 2025) employed
 these metrics to analyse multi-head attention mechanisms in
 transformers. Surprisingly, such frameworks are largely un-
 studied in mathematics, aside from expander graphs (Csóka
 & Grabowski, 2022) and the spectral analysis of Stanković
 et al. (2025). (Herasimchyk et al., 2026) and (Chowdhury,
 2026) used an *influence* metric to analyse the Lost-In-The-
 Middle effect in transformers, which is closely related to
 the normalised minimax fidelity (Barbero et al., 2024).

3. Background

To formalise the study of architecture design in sequence
 processing, we model the forward pass of a neural network
 as a directed acyclic graph (DAG) augmented with self-
 loops, referred to as a feedforward computational graph
 $G = (\mathcal{V}, \mathcal{E})$. The nodes \mathcal{V} store the feature representations
 at various layers, and the directed edges \mathcal{E} dictate the flow
 of computation. Further, we assume that the nodes in \mathcal{V} are
 topologically sorted such that $\mathcal{V} = \{1, \dots, |\mathcal{V}|\}$ where node
 1 is the unique source and node $|\mathcal{V}|$ is the unique sink n . If a
 directed edge exists from i to j , then $i \leq j$. The adjacency
 matrix \mathbf{A} is defined in the usual way: $a_{ij} = 1 \Leftrightarrow (j, i) \in \mathcal{E}$
 and $a_{ij} = 0$ otherwise.

Following (Vitvitskyi et al., 2025), our analysis relies on
 certain structural assumptions on G :

- **Self-Loops.** Every node possesses a self-loop, i.e.,
 $(i, i) \in \mathcal{E}$ for all $i \in \mathcal{V}$, representing the property that a
 node can preserve its own state across computational
 steps.

- **Unique Sink.** There exists a unique sink node, denoted as $n \in V$, which is the only node with an out-degree of one (its self-loop). This represents the final output representation of the network.
- **Unique Source.** Unlike (Vitvitskyi et al., 2025), we omit the self-similarity assumption and replace it with a requirement that the graph has a defined starting point. We assume there exists a unique *source*, representing a unified starting point for the flow of information. Consequently, due to the presence of self-loops, the source is the only node with in-degree one.

Below, we recap the (averaged) mixing time and (normalised) minimax fidelity framework introduced by (Vitvitskyi et al., 2025).

3.1. Quantifying Information Propagation Speed: Mixing Time

Let $\mathbf{W} \in \mathbb{R}^{n \times n}$ be the column-normalised forward walk matrix, where w_{uv} is the u -th row and v -th column entry of \mathbf{W} and represents the probability that a random walker at node v transitions to node u . Similarly, w_{uv}^t represents that a random walker at node v transitions to node u in t steps. The forward mixing time $T_{fwd}(G)$ is defined as the minimum number of steps t from *any* starting probability distribution \mathbf{x} to converge sufficiently close to the stationary distribution $\boldsymbol{\pi}$, which in our case is one-hot in the unique sink n :

$$T_{fwd}(G) = \min \left\{ t \geq 0 : \max_{\mathbf{x}} [\|\mathbf{W}^t \mathbf{x} - \boldsymbol{\pi}\|_1] < \frac{1}{2} \right\}$$

We note that while (Vitvitskyi et al., 2025) analyse *averaged* mixing times (Espuny Díaz et al., 2024)—averaging the minimum steps t required to reach the stationary distribution across all starting nodes—our analysis adopts a *worst-case* analysis. This worst-case formulation (represented by the maximisation over \mathbf{x}) establishes higher parity with standard Markov chain literature.

3.2. Quantifying Information Sharpness: Normalised Minimax Fidelity

Let $\Delta \in \mathbb{R}^{n \times n}$ be the row-normalised diffusion matrix, where Δ_{ij} corresponds to the proportion of influence a prior node j exerts on node i during a single step of information mixing. Over t timesteps, the fidelity for a node i , Δ_{ni}^t , quantifies the total accumulated portion of information from a source node i that successfully reaches the sink n .

The maximal fidelity for a given node i , defined as $\phi_i = \max_t \Delta_{ni}^t$, represents the highest concentration of its features captured at the sink over the entire duration of the diffusion. To capture the robustness bottleneck of the entire

graph, the minimax fidelity computes the minimum of these maximum fidelities across all nodes, normalised by the size of the graph n :

$$F(G) = n \cdot \min_{i \in V} \max_t \Delta_{ni}^t$$

Intuitively, this metric identifies the “weakest link” in the graph: the node whose best-case representation at the sink is the most attenuated.

3.3. Edge Weighting Schemes: Uniform vs. Non-Uniform Graphs

Having established the general metrics for evaluating a computational graph, we distinguish between two classes of graphs based on how nodes aggregate incoming information. This categorisation defines how the diffusion matrix Δ and the random walk matrix \mathbf{W} are populated.

- **Uniform Graphs:** Every node assigns equal weight to all of its incoming edges. Let $d_{\rightarrow i}$ denote the in-degree of node i , and let $d_{j \rightarrow}$ denote the out-degree of node j . The diffusion matrix uniformly distributes incoming influence, while the forward walk matrix distributes outgoing transition probabilities equally. Formally, for any edge $(j, i) \in \mathcal{E}$:

$$\Delta_{ij} = \begin{cases} 1/d_{\rightarrow i} & \text{if } (j, i) \in \mathcal{E} \\ 0 & \text{otherwise} \end{cases}$$

$$w_{ij} = \begin{cases} 1/d_{j \rightarrow} & \text{if } (j, i) \in \mathcal{E} \\ 0 & \text{otherwise} \end{cases}$$

- **Non-Uniform Graphs:** Nodes can assign arbitrary, data-dependent weights to their incoming connections (e.g., standard attention mechanisms). In this setting, the diffusion matrix Δ can be any arbitrary row-stochastic matrix that respects the graph’s topology (i.e., $\Delta_{ij} \geq 0$, and $\sum_j \Delta_{ij} = 1$). To maintain a consistent theoretical duality between backward information flow and the forward random walk, in this setting, we define the non-uniform forward walk matrix \mathbf{W} as the column-normalisation of the diffusion matrix Δ :

$$w_{ij} = \begin{cases} \Delta_{ij}/C_j & \text{if } (j, i) \in \mathcal{E} \\ 0 & \text{otherwise} \end{cases}$$

where $C_j = \sum_k \Delta_{kj}$ acts as the normalising constant for the outgoing transitions from node j .

Assumption: c -constrained non-uniform graphs. In line with what was done in (Herasimchyk et al., 2026), we make the assumption that non-uniform graphs are c -constrained. We define a graph as c -constrained if every valid edge weight

in the diffusion matrix satisfies $\Delta_{ij} > c$ for some small constant $c > 0$. In other words, we bound the graph away from zero by c . This assumption is justified, as most edge weights are derived from attention mechanisms that employ a `softmax` function, and therefore no edge weight can ever be strictly zero.

4. Equivalence of Fidelity and the ‘Backwards’ Random Walk

Before presenting our main theoretical bounds, we first establish an equivalence between information diffusion in a feedforward graph and a random walk on its reversed counterpart. Importantly, the results in this section hold for both uniform graphs and non-uniform graphs with arbitrary edge weightings.

Definition 4.1. Let $G = (\mathcal{V}, \mathcal{E})$ be a computational feedforward graph. Define $\tilde{G} = (\mathcal{V}, \tilde{\mathcal{E}})$ to be its reversed counterpart, where nodes are relabeled such that node i in G maps to node $n - i + 1$ in \tilde{G} . The edges are correspondingly reversed: $\forall i, j : (i, j) \in \mathcal{E} \iff (n - j + 1, n - i + 1) \in \tilde{\mathcal{E}}$. Let $\tilde{\mathbf{W}}$ be the column-normalised random walk matrix of its reversed counterpart \tilde{G} and $r(i) = n - i + 1$. Then the edge weights are defined such that, for all $(j, i) \in \mathcal{E}$, we have $\Delta_{ij} = \tilde{w}_{r(j), r(i)}$ and $\tilde{\Delta}_{ij} = \tilde{w}_{ij} / \tilde{C}_i$ where $\tilde{C}_i = \sum_k \tilde{w}_{ik}$.

By reversing the graph, the unique forward sink of G (node n) becomes the source node of \tilde{G} , and the unique source of G (node 1) becomes the sink node of \tilde{G} . We now show that computing the fidelity of G follows directly from computing a random walk on \tilde{G} .

Theorem 4.2. For any number of steps t and any nodes i, j , we have:

$$\Delta_{ij}^t = \tilde{w}_{r(j), r(i)}^t$$

Proof Sketch. The proof proceeds by straightforward induction on t . The base case ($t = 1$) holds by definition. For the inductive step, we delegate the details to Section A. \square

Given this, we can arrive at the following straightforward corollary:

Corollary 4.3. The last row of Δ of any computational feedforward G evolves exactly like the first column of $\tilde{\mathbf{W}}$.

5. The Backwards Mixing Time - Fidelity Bound

Having established the duality between the forward diffusion of information and the backward random walk, we now present our main theoretical results. We first prove the central result of this paper: the backwards mixing time forms an upper bound over normalised minimax fidelity. We then show that this bound is tight for a certain family of graphs we call *conjoined star* graphs.

5.1. Fidelity is Bounded by Backward Mixing Time

We utilise our equivalence between the fidelity and the backward random walk proven in the earlier section (Theorem 4.2) to show that the backwards mixing time forms an upper bound on the fidelity:

Theorem 5.1 (Fidelity-Backwards-Mixing Bound). Let G be a feedforward graph with a unique source and a unique sink with $n \geq 2$. Let $F(G)$ be its minimax fidelity and $T_{bwd}(G)$ be the worst-case mixing time of its backward random walk. Then, for both uniform and non-uniform graphs:

$$F(G) \leq \frac{8}{3} \cdot T_{bwd}(G)$$

Proof. We first recall the definition of the worst-case mixing time $T_{bwd}(G)$ to be the lowest t such that for any starting distribution \mathbf{x} and stationary distribution $\boldsymbol{\pi}$:

$$\|\tilde{\mathbf{W}}^t \mathbf{x} - \boldsymbol{\pi}\|_1 < \frac{1}{2}$$

where $\tilde{\mathbf{W}}$ is the transition matrix of the reverse random walk. Because there is a unique forward sink n (which becomes the unique source node in the reverse graph), the stationary distribution $\boldsymbol{\pi}$ is exactly 1 at the absorbing node and 0 everywhere else.

By the definition of mixing time, after $T_{bwd}(G)$ steps, the probability of the walker being at any node other than the absorbing node is at most $1/4$. More generally, after $k \cdot T_{bwd}(G)$ steps, this probability drops to at most $(1/4)^k$.

In our case, since there is a unique source and sink, the stationary distribution $\boldsymbol{\pi}$ for both the forward and reverse random walk is 1 in the absorbing node and 0 everywhere else. Let S be the absorbing set (a singleton) in the reverse walk, and let $P(X_t \notin S)$ be the probability that the walker has not reached the absorbing node after t steps. Recall that $\phi_i = \max_t \Delta_{ni}^t$. Because the maximum of any non-negative sequence is bounded by its sum ($\phi_i \leq \sum_{t=0}^{\infty} \Delta_{ni}^t$), we can bound the sum of fidelities for all non-sink nodes in

\tilde{G} :

$$\begin{aligned}
 \sum_{i \notin S} \phi_i &\leq \sum_{i \notin S} \sum_{t=0}^{\infty} \Delta_{ni}^t \\
 &= \sum_{t=0}^{\infty} \sum_{i \notin S} \Delta_{ni}^t \\
 &= \sum_{t=0}^{\infty} \sum_{i \notin S} \tilde{w}_{r(i), r(n)}^t \\
 &= \sum_{t=0}^{\infty} P(X_t \notin S) \\
 &= \sum_{k=0}^{\infty} \sum_{r=0}^{T_{bwd}(G)-1} P(X_{k \cdot T_{bwd}(G)+r} \notin S)
 \end{aligned}$$

Because all backward random walks must eventually culminate at the unique sink, $P(X_t \notin S)$ is monotonically decreasing with t . Thus, $P(X_t \notin S) \geq P(X_{t+1} \notin S)$, allowing us to upper-bound the value of each block of size $T_{bwd}(G)$ by its starting value:

$$\begin{aligned}
 \sum_{i \notin S} \phi_i &\leq \sum_{k=0}^{\infty} T_{bwd}(G) \cdot P(X_{k \cdot T_{bwd}(G)} \notin S) \\
 &< \sum_{k=0}^{\infty} T_{bwd}(G) \cdot \left(\frac{1}{4}\right)^k \\
 &= T_{bwd}(G) \cdot \frac{1}{1 - 1/4} = \frac{4}{3} \cdot T_{bwd}(G)
 \end{aligned}$$

Putting it all together (by bounding the minimum of the maximum node fidelities by its average), we can bound the minimax fidelity $F(G)$:

$$\begin{aligned}
 F(G) &= n \cdot \min_{i \in V} \max_t \Delta_{ni}^t = n \cdot \min_{i \in V} \phi_i \\
 &\leq n \cdot \min_{i \notin S} \phi_i \\
 &\leq n \cdot \frac{1}{n-1} \sum_{i \notin S} \phi_i \\
 &\leq n \cdot \frac{1}{n-1} \cdot \frac{4}{3} \cdot T_{bwd}(G) \\
 &\leq \frac{8}{3} \cdot T_{bwd}(G)
 \end{aligned}$$

where the final inequality holds because $\frac{n}{n-1} \leq 2$ for $n \geq 2$, and thus $2 \cdot \frac{4}{3} = \frac{8}{3}$ for any graph where $n \geq 2$. \square

5.2. Tightness of the Tradeoff

We demonstrate that the above bound is tight. To establish this, it is sufficient to identify a class of uniform graphs where $F(G)$ and $T_{fwd}(G)$ are equivalent up to constants. Because uniform graphs represent a strict subset of non-uniform graphs (specifically, the case where all edge weights

are constrained to be equal), proving tightness in the uniform setting proves tightness for the non-uniform setting too – guaranteeing that our tradeoff bounds cannot be improved further.

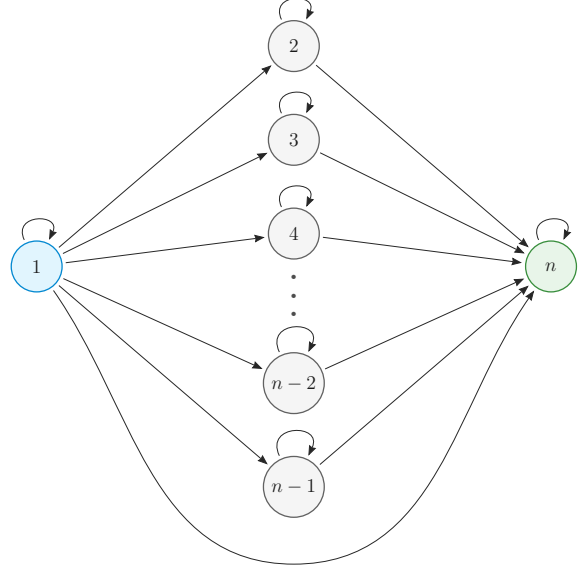


Figure 2. The conjoined star graph.

Theorem 5.2. *The Fidelity-Backward-Mixing Bound is tight up to asymptotic constants.*

Proof. Define the class of *conjoined star* graphs as uniform feedforward graphs where the unique source (node 1) is connected to every other node, and every other node is connected to the unique sink (node n). Intuitively, this is a combination of two star graphs: an *outward* star centred at node 1, and an *inward* star centred at node n (Figure 2). Formally, assuming self-loops, the edge set \mathcal{E} is defined as:

$$\begin{aligned}
 \mathcal{E} &= \{(1, i) \mid i \in (1, n]\} \\
 &\cup \{(i, n) \mid i \in [1, n)\} \\
 &\cup \{(i, i) \mid i \in [1, n]\}
 \end{aligned}$$

Because this graph is uniform, it is easy to see that $T_{fwd}(G) = T_{bwd}(G)$. For convenience, we will analyse the forward mixing time. To bound $T_{fwd}(G)$, we note that the random walk would be slowest if it started at node 1, the source node. It has a probability of $\frac{n-1}{n}$ of leaving the source node at a single step, and hence will take $\frac{n}{n-1}$ steps on average to leave the source node. Once it leaves the source node, then any node except the source node will take at most 2 steps in expectation to reach the sink, as all intermediary nodes have $d_{out} = 2$. Since $\frac{n}{n-1} \leq 2$ for $n \geq 2$, we have $T_{fwd} = O(1)$. Using Theorem 4.2, we note that the fidelity of a node u evolves exactly like the first column

of the reverse random walk matrix. Hence, the maximum fidelity ϕ_u is the maximum probability the random walker starting from n is at node u during its course of the random walk. The maximum probability for the random walker to be at the node u is at the first step, which is $\frac{1}{n}$. This means that $\phi_u = \frac{1}{n}$ for all intermediary nodes u , ($1 < u < n$), and normalised minimax fidelity is equal to 1, which implies $F(G) = \Theta(1) = T_{fwd}(G) = T_{bwd}(G)$. \square

6. Relating Backwards and Forwards Mixing Time

Despite having a simple proof, we note that the above result is surprising, as it is not that obvious that the backwards mixing time can be connected to fidelity at first glance.

6.1. A Fundamental Tradeoff for Symmetric Graphs

Using this bound, it follows that a fundamental tradeoff bound between the forwards mixing time and the normalised minimax fidelity exists for *symmetric* graphs in the uniform setting, where $T_{fwd}(G) = T_{bwd}(G)$.

In the case where G is isomorphic to \tilde{G} , i.e. G is equivalent to the reversed graph \tilde{G} up to node relabeling, we would have an equivalence between the backwards mixing time and forward mixing time ($T_{bwd}(G) = T_{fwd}(G)$) and we deem that the graph is *symmetric*. In the uniform setting, this occurs for a wide range of attentional masks proposed in the literature, i.e. fully-connected attention (Vaswani et al., 2023), sliding window attention (Wang et al., 2019), LogSparse (Li et al., 2020), strided attention (Child et al., 2019), and FiboAttention (Rahimian et al., 2026). We also note that the FunSearch graph proposed by (Vitvitskyi et al., 2025) is symmetric. When the graph is symmetric, we will have $F(G) \leq \frac{8}{3} \cdot T_{fwd}(G)$. Hence, for symmetric graphs, this demonstrates a fundamental tradeoff: a computational graph cannot simultaneously achieve optimal forward mixing speed and maximal representation sharpness.

6.2. General Tradeoffs

More generally, for non-symmetric uniform graphs, e.g. LongFormer (Beltagy et al., 2020), BigBird (Zaheer et al., 2021), and non-uniform graphs, if $T_{bwd}(G)$ is related to $T_{fwd}(G)$ by some f , then we arrive at the following general corollary:

Corollary 6.1 (General Fidelity-Forward-Mixing Bound.). *If the backward mixing time can be bounded by the forward mixing time via some function f , such that $T_{bwd}(G) = f(T_{fwd}(G))$, then the minimax fidelity is bottlenecked by the forward mixing time:*

$$F(G) \leq \frac{8}{3} \cdot f(T_{fwd}(G))$$

This bound assumes that an f that relates T_{fwd} and T_{bwd} must exist. For completeness, we next demonstrate that such a bounding function f must exist for both uniform and non-uniform graphs through a simple argument that upper-bounds both $T_{fwd}(G)$ and $T_{bwd}(G)$.

Remark. From the definition of fidelity and from Theorem 4.2, we can see that $\phi_i \leq 1$ for every node i , as the probability to be at node i at the backward random walk is at most 1 at any timestep t . Hence, we can arrive at a trivial bound for normalised minimax fidelity: $F(G) \leq n$ for all G . We acknowledge that the bounds that are presented below imply a looser bound for $F(G)$ than the trivial $O(n)$ bound. The purpose of including these bounds is to show that a bounding function f relating forward and backward must exist and does not blow up.

6.3. Mixing Time Bounds for Uniform Graphs

For uniform graphs, we find that the forward and backward mixing times are polynomially equivalent, scaling at most linearly with each other with the sequence length n .

Theorem 6.2. *For any uniform feedforward graph with $n \geq 2$ nodes, self-loops, and unique source and sink, the forward and backward mixing times linearly bound one another:*

$$\begin{aligned} T_{bwd}(G) &\leq 8n \cdot T_{fwd}(G) \\ T_{fwd}(G) &\leq 8n \cdot T_{bwd}(G) \end{aligned}$$

Proof Sketch. Because nodes in uniform graphs weight all incoming edges equally, the maximal probability of a self-loop is $1/2$. This allows us to model leaving any state as a geometric random variable. Applying Markov’s inequality across all n nodes yields an absolute maximum bound of $T \leq 8n$. The full derivation is provided in Section B. \square

6.4. Mixing Time Bounds for Non-Uniform Graphs

We now extend the above results for non-uniform graphs:

Theorem 6.3. *For any c -constrained non-uniform feedforward graph with $n \geq 2$, we have*

$$\begin{aligned} T_{bwd}(G) &\leq \frac{4n}{c} \cdot T_{fwd}(G) \\ T_{fwd}(G) &\leq \frac{4n}{c} \cdot T_{bwd}(G) \end{aligned}$$

Proof Sketch. We analyse the expected hitting times of the random walks. Because the diffusion weights are lower-bounded by c , the sum of fractional coefficients at any node acts as a convex combination bounded by the maximum previous term plus $1/c$. Using Markov’s inequality on these

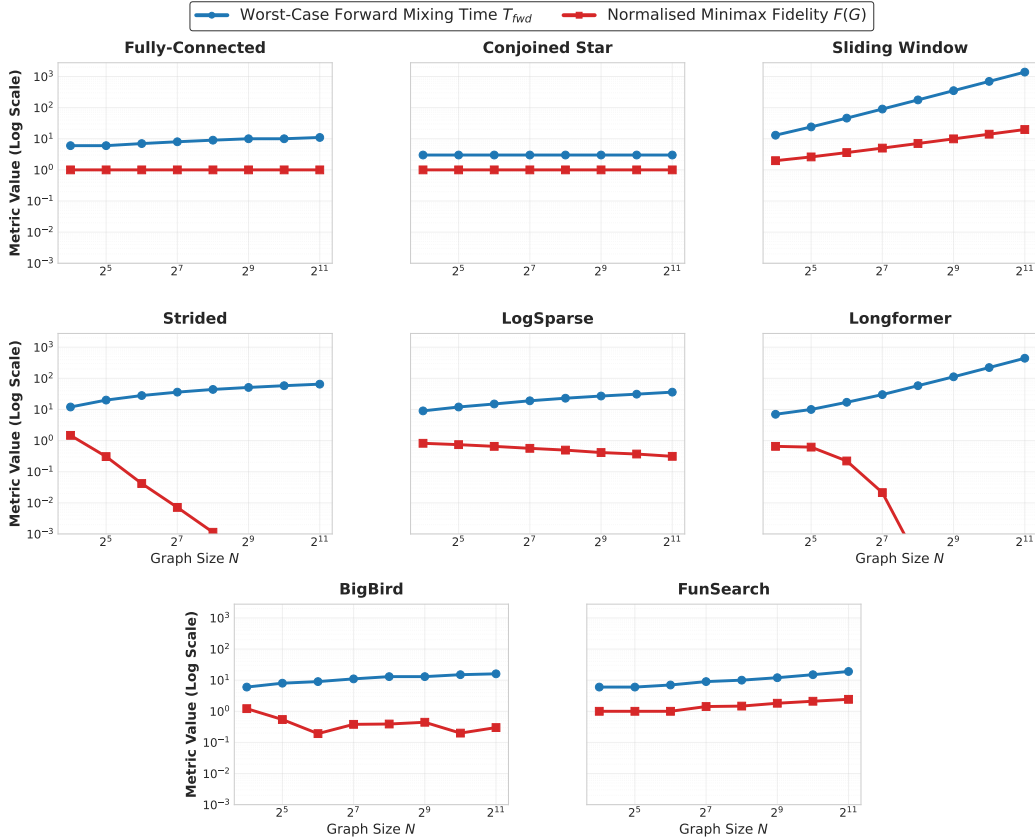


Figure 3. Mixing time and normalised minimax fidelity for a range of attentional masks.

partial sums yields $T \leq 4n/c$. The complete proof is in Section C. \square

Remark on the Bounding Function f . Using Theorem 6.2 and Theorem 6.3 and plugging it to Theorem 6.1, we can arrive at general upper bounds for fidelity: $F(G) \leq \frac{64}{3} \cdot n \cdot T_{fwd}(G)$ for uniform graphs and $F(G) \leq \frac{32n}{3c} \cdot T_{fwd}(G)$ for non-uniform graphs. It is crucial to emphasise that these bounds do not serve as the “final say” on the tightness of the tradeoff for all specific architectures.

For specific, well-structured classes of computational graphs, the exact bounding function f can be significantly tighter. For example, as discussed earlier, in any symmetric computational graph, $T_{fwd}(G) = T_{bwd}(G)$, making the bounding function the simple identity $f(x) = x$. In such cases, the tradeoff collapses to a drastically tighter, constant-factor bound: $F(G) \leq \frac{8}{3} \cdot T_{fwd}(G)$. Thus, the general bounds above establish a worst-case guarantee, while highlighting that the severity of the tradeoff is dictated by the specific structure of the chosen feedforward graph.

7. Empirical Demonstration: Uniform Graphs

To demonstrate our bounds hold, in Figure 3, we plot the worst-case forward mixing time $T_{fwd}(G)$ and the minimax fidelity $F(G)$ for a range of attentional masks in the uniform graph setting. Specifically, we plot symmetric masks—fully-connected attention, conjoined star, sliding window attention, strided attention, LogSparse attention, and the FunSearch mask—and non-symmetric masks—LongFormer and BigBird. Each point in the graph corresponds to either minimax fidelity or worst-case mixing time, with varying graph sizes $N \in \{16, 32, 64, \dots, 2048\}$. We use a default window size of 3 for masks with a sliding window component. For strided attention, we use a stride size of 8, and for LongFormer we use a dilation size of 3. For BigBird, we add 5 random edges to each node.

From Figure 3, we can see for symmetric graphs, it is evident that forward mixing time forms an upper bound on the normalised minimax fidelity, because $T_{fwd}(G) = T_{bwd}(G)$.

For non-symmetric graphs, we also observe the same effect due to Theorem 6.2. We also note that in practice, the gap between normalised minimax fidelity and forward mixing time is significantly tighter than the bound in Theorem 6.2. We defer the plotting of these masks in the non-uniform setting for future work, as it is unclear how to aggregate the large number of possible attentional masks in this setting into a single plot.

8. Conclusion

In this work, we establish an equivalence between the forward diffusion of representations and a random walk on the reversed graph. Using this, we proved that the minimax fidelity of a computational graph is tightly upper-bounded by its backward mixing time, $F(G) \leq \frac{8}{3} \cdot T_{bwd}(G)$ —establishing the backwards mixing time as a metric to analyse the sharpness of information propagation.

Next, for symmetric graphs in the uniform setting, we demonstrate a fundamental tradeoff: a computational graph cannot simultaneously achieve optimal forward mixing speed and maximal representation sharpness. In the more general case, we proved that $F(G) \leq f(T_{fwd}(G))$ where $f(T_{fwd}(G))$ is bounded by $O(n) \cdot T_{fwd}(G)$ in the uniform setting and $O(n/c) \cdot T_{fwd}(G)$ in the non-uniform setting.

We note limitations that present clear avenues for future work:

- **Worst-Case vs. Average-Case.** To provide absolute mathematical guarantees on minimax fidelity, our proofs necessitate the use of worst-case mixing times. Because sequence models often operate over distributions where average-case performance is sufficient, a highly relevant extension would be to formally define an *averaged* fidelity metric and investigate whether a parallel tradeoff exists between averaged fidelity and averaged mixing time.
- **Analysing Graphs With Multiple Forward Sources.** Our results assume that feedforward graphs have a unique source. Although this captures most masks in practice, some masks have multiple sources, e.g. the blockwise attention in BlockBERT (Qiu et al., 2020). A natural next step would be to extend our analysis for such graphs.
- **Deriving Tighter Universal Bounds for Mixing Times.** While we proved that the bounding function relating forward and backward mixing times scales at most by $O(n)$ and $O(n/c)$, these are very loose universal bounds. Future work could incorporate additional structural assumptions—such as bounded graph width, specific sparsity patterns, or expansion properties—to

derive a tighter universal bound that more closely realises practical constraints and demands.

Impact Statement

This paper proves the central result that normalised minimax fidelity is upper-bounded by backwards mixing time. In isolation, this work does not present any direct societal risks or negative ethical consequences. However, we predict that our work will inform researchers of this inherent property in computational feedforward graphs and to inspire the design of better models. Hence, any societal risks that our work poses should be attributed to the consequences of accelerating the design of such models.

References

- Attali, H., Buscaldi, D., and Pernelle, N. Rewiring techniques to mitigate oversquashing and oversmoothing in gnn: A survey, 2024. URL <https://arxiv.org/abs/2411.17429>.
- Barbero, F., Banino, A., Kapturowski, S., Kumaran, D., Araújo, J. G., Vitvitskyi, A., Pascanu, R., and Veličković, P. Transformers need glasses! information oversquashing in language tasks. *Advances in Neural Information Processing Systems*, 37:98111–98142, 2024.
- Beltagy, I., Peters, M. E., and Cohan, A. Longformer: The long-document transformer, 2020. URL <https://arxiv.org/abs/2004.05150>.
- Chen, L., Xu, D., An, C., Wang, X., Zhang, Y., Chen, J., Liang, Z., Wei, F., Liang, J., Xiao, Y., and Wang, W. Powerattention: Exponentially scaling of receptive fields for effective sparse attention, 2025. URL <https://arxiv.org/abs/2503.03588>.
- Child, R., Gray, S., Radford, A., and Sutskever, I. Generating long sequences with sparse transformers, 2019. URL <https://arxiv.org/abs/1904.10509>.
- Chowdhury, B. D. Lost in the middle at birth: An exact theory of transformer position bias, 2026. URL <https://arxiv.org/abs/2603.10123>.
- Csóka, E. and Grabowski, Ł. On directed analogues of expander and hyperfinite graph sequences. *Combinatorics, Probability and Computing*, 31(2):184–197, 2022.
- de Ocaíz Borde, H. S. Beyond parallelism: Synergistic computational graph effects in multi-head attention, 2025. URL <https://arxiv.org/abs/2507.02944>.
- Deac, A., Lackenby, M., and Veličković, P. Expander graph propagation. In *Learning on Graphs Conference*, pp. 38–1. PMLR, 2022.

- 440 Espuny Díaz, A., Morris, P., Perarnau, G., and Serra, O.
441 Speeding up random walk mixing by starting from a
442 uniform vertex. *Electronic journal of probability*, 29:
443 1–25, 2024.
- 444
445 Giovanni, F. D., Giusti, L., Barbero, F., Luise, G., Lio',
446 P., and Bronstein, M. On over-squashing in message
447 passing neural networks: The impact of width, depth, and
448 topology, 2023. URL [https://arxiv.org/abs/
449 2302.02941](https://arxiv.org/abs/2302.02941).
- 450
451 Herasimchyk, H., Labryga, R., Prusina, T., and Laue, S.
452 A residual-aware theory of position bias in transform-
453 ers, 2026. URL [https://arxiv.org/abs/2602.
454 16837](https://arxiv.org/abs/2602.16837).
- 455
456 Huang, S., Poursafaei, F., Danovitch, J., Fey, M., Hu, W.,
457 Rossi, E., Leskovec, J., Bronstein, M., Rabusseau, G.,
458 and Rabbany, R. Temporal graph benchmark for machine
459 learning on temporal graphs, 2023. URL [https://
460 arxiv.org/abs/2307.01026](https://arxiv.org/abs/2307.01026).
- 461
462 Kim, J., Kim, H., Kim, H., Lee, D., and Yoon, S. A compre-
463 hensive survey of deep learning for time series forecast-
464 ing: Architectural diversity and open challenges, 2025.
465 URL <https://arxiv.org/abs/2411.05793>.
- 466
467 Li, S., Jin, X., Xuan, Y., Zhou, X., Chen, W., Wang, Y.-
468 X., and Yan, X. Enhancing the locality and breaking the
469 memory bottleneck of transformer on time series forecast-
470 ing, 2020. URL [https://arxiv.org/abs/1907.
471 00235](https://arxiv.org/abs/1907.00235).
- 472
473 Qiu, J., Ma, H., Levy, O., tau Yih, S. W., Wang, S., and Tang,
474 J. Blockwise self-attention for long document understand-
475 ing, 2020. URL [https://arxiv.org/abs/1911.
476 02972](https://arxiv.org/abs/1911.02972).
- 477
478 Rahimian, A. K., Govind, M. K., Maity, S., Reilly, D.,
479 Kümmerle, C., Das, S., and Dutta, A. Fibottention: In-
480 ceptive visual representation learning with diverse atten-
481 tion across heads, 2026. URL [https://arxiv.org/
482 abs/2406.19391](https://arxiv.org/abs/2406.19391).
- 483
484 Stanković, L., Daković, M., Bardi, A. B., Brajović, M.,
485 and Stanković, I. Fourier analysis of signals on directed
486 acyclic graphs (dag) using graph zero-padding. *Digital
487 Signal Processing*, 159:104995, 2025. ISSN 1051-
488 2004. doi: <https://doi.org/10.1016/j.dsp.2025.104995>.
489 URL [https://www.sciencedirect.com/
490 science/article/pii/S105120042500017X](https://www.sciencedirect.com/science/article/pii/S105120042500017X).
- 491
492 Sutskever, I., Vinyals, O., and Le, Q. V. Sequence to
493 sequence learning with neural networks, 2014. URL
494 <https://arxiv.org/abs/1409.3215>.
- Topping, J., Giovanni, F. D., Chamberlain, B. P., Dong,
X., and Bronstein, M. M. Understanding over-squashing
and bottlenecks on graphs via curvature, 2022. URL
<https://arxiv.org/abs/2111.14522>.
- van den Oord, A., Dieleman, S., Zen, H., Simonyan, K.,
Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A.,
and Kavukcuoglu, K. Wavenet: A generative model for
raw audio, 2016. URL [https://arxiv.org/abs/
1609.03499](https://arxiv.org/abs/1609.03499).
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones,
L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention
is all you need, 2023. URL [https://arxiv.org/
abs/1706.03762](https://arxiv.org/abs/1706.03762).
- Vitvitskyi, A., Araújo, J. G. M., Lackenby, M., and
Veličković, P. What makes a good feedforward compu-
tational graph?, 2025. URL [https://arxiv.org/
abs/2502.06751](https://arxiv.org/abs/2502.06751).
- Wang, Z., Ng, P., Ma, X., Nallapati, R., and Xiang, B.
Multi-passage bert: A globally normalized bert model for
open-domain question answering, 2019. URL [https://
arxiv.org/abs/1908.08167](https://arxiv.org/abs/1908.08167).
- Wiedemer, T., Li, Y., Vicol, P., Gu, S. S., Matarese, N.,
Swersky, K., Kim, B., Jaini, P., and Geirhos, R. Video
models are zero-shot learners and reasoners, 2025. URL
<https://arxiv.org/abs/2509.20328>.
- Wilson, J., Bechler-Speicher, M., and Veličković, P. Cayley
graph propagation. *arXiv preprint arXiv:2410.03424*,
2024.
- Xia, Z., Pan, X., Song, S., Li, L. E., and Huang, G. Vi-
sion transformer with deformable attention, 2022. URL
<https://arxiv.org/abs/2201.00520>.
- Zaheer, M., Guruganesh, G., Dubey, A., Ainslie, J., Albeti,
C., Ontanon, S., Pham, P., Ravula, A., Wang, Q., Yang,
L., and Ahmed, A. Big bird: Transformers for longer
sequences, 2021. URL [https://arxiv.org/abs/
2007.14062](https://arxiv.org/abs/2007.14062).

A. Proof of Theorem 4.2

Theorem 4.2. *Let Δ be the row-stochastic diffusion matrix of an arbitrary feedforward graph G , let $\tilde{\mathbf{W}}$ be the column-normalized random walk matrix of its reversed counterpart \tilde{G} , and let $r(i) = |V| - i + 1$. For any number of steps t and any nodes i, j , we have $\Delta_{i,j}^t = \tilde{w}_{r(j),r(i)}^t$.*

Proof. We proceed by induction on t .

The base case holds by definition. For the inductive step, assume that the theorem holds for some arbitrary step $t \geq 1$. That is, our inductive hypothesis states that for all $i, j \in \mathcal{V}$:

$$\Delta_{i,j}^t = \tilde{w}_{r(j),r(i)}^t$$

We must show this implies $\Delta_{i,j}^{t+1} = \tilde{w}_{r(j),r(i)}^{t+1}$.

Using the standard definition of matrix multiplication (and the property that matrix powers commute), we can expand $\Delta_{i,j}^{t+1}$ by summing over all possible intermediate nodes k :

$$\Delta_{i,j}^{t+1} = \sum_{k \in V} \Delta_{i,k}^t \cdot \Delta_{k,j} \quad (1)$$

Similarly, we expand $\tilde{w}_{r(j),r(i)}^{t+1}$ by summing over all possible intermediate nodes $m \in \mathcal{V}$:

$$\tilde{w}_{r(j),r(i)}^{t+1} = \sum_{m \in V} \tilde{w}_{r(j),m} \cdot \tilde{w}_{m,r(i)}^t$$

which is equivalent to:

$$\tilde{w}_{r(j),r(i)}^{t+1} = \sum_{k \in V} \tilde{w}_{r(j),r(k)} \cdot \tilde{w}_{r(k),r(i)}^t \quad (2)$$

We can now apply our previously established equivalencies to Equation 2. By the base case, we know $\tilde{w}_{r(j),r(k)} = \Delta_{k,j}$. By the inductive hypothesis, we know $\tilde{w}_{r(k),r(i)}^t = \Delta_{i,k}^t$. Substituting these into the summation gives:

$$\begin{aligned} \tilde{w}_{r(j),r(i)}^{t+1} &= \sum_{k \in V} \Delta_{k,j} \cdot \Delta_{i,k}^t \\ &= \sum_{k \in V} \Delta_{i,k}^t \cdot \Delta_{k,j} \\ &= \Delta_{i,j}^{t+1} \end{aligned}$$

B. Proof of Theorem 6.2

To prove Theorem 6.2, we first establish global upper and lower bounds on the mixing times for uniform graphs.

Lemma B.1. *For any uniform feedforward graph with self-loops and unique source and sink, the mixing times satisfy $T_{fwd}(G) \leq 8n$ and $T_{bwd}(G) \leq 8n$.*

Proof. Because the graph is uniform, a node weights all incoming edges equally. The presence of a self-loop means the maximal probability of remaining at the current state is $1/2$, as every non-sink node has an out-degree of at least 2 (a self-loop plus one outgoing edge). Thus, the probability of leaving any state is at least $1/2$. Consequently, leaving a state can be modeled as a geometric random variable, meaning the expected number of steps to leave a node is at most 2. Summing this expectation across a maximum path length of n nodes, the expected number of steps to reach the target node is $\mathbb{E}[T] \leq 2n$. By Markov's inequality, the probability of failing to reach the sink after $8n$ steps is bounded by $\mathbb{P}(T \geq 8n) \leq \frac{2n}{8n} = \frac{1}{4}$. By the definition of mixing time, $T_{fwd} \leq 8n$. This logic holds identically for both the forward and backward random walks. \square

Proof of Theorem 6.2. By Lemma B.1, $T_{bwd} \leq 8n$. Since $n \geq 2$, we also have $T_{fwd} \geq 1$ and $T_{bwd} = 1$. Substituting this inequality yields:

$$T_{bwd} \leq 8n \cdot 1 \leq 8n \cdot T_{fwd}$$

The derivation for bounding T_{fwd} by T_{bwd} follows identically. \square

C. Proof of Theorem 6.3

We define a graph as c -constrained if for any $(i, j) \in \mathcal{E}$, then $\Delta_{j,i} > c$ for some constant $c > 0$. We first bound the absolute expected hitting times for these graphs.

Lemma C.1. *For any c -constrained non-uniform graph, $T_{bwd} \leq \frac{4n}{c}$ and $T_{fwd} \leq \frac{4n}{c}$.*

Proof. Let $h_i = \mathbb{E}[T_{bwd} \mid X_0 = i]$ be the expected time for the backward walk to reach the source, with $h_1 = 0$. By standard Markov chain properties:

$$\begin{aligned} h_i &= 1 + \Delta_{i,i} h_i + \sum_{j < i} \Delta_{i,j} h_j \\ (1 - \Delta_{i,i}) h_i &= 1 + \sum_{j < i} \Delta_{i,j} h_j \\ h_i &= \frac{1}{\sum_{j < i} \Delta_{i,j}} + \sum_{j < i} \frac{\Delta_{i,j}}{\sum_{k < i} \Delta_{i,k}} h_j \end{aligned}$$

Because node i is not the source, it must have at least one backward edge to some $j < i$. By the c -constraint, $\sum_{j < i} \Delta_{i,j} \geq c$. Because the summation of fractional coefficients strictly equals 1, the second term is a convex combination. Thus, it is bounded by the maximum previous term:

$$h_i \leq \frac{1}{c} + \max_{j < i} h_j$$

By unrolling the recursion, we have $\max_i h_i \leq \frac{n}{c}$. By Markov's Inequality, enforcing $\mathbb{P}(T > t) \leq \frac{\mathbb{E}[T]}{t} \leq 1/4$ requires $t \geq \frac{4n}{c}$, and hence $T_{bwd} \leq \frac{4n}{c}$.

For the forward walk, let $g_j = \mathbb{E}[T_{fwd} \mid X_0 = j]$, with $g_n = 0$. Denote $w_{j,k}$ for the j -th row and k -th column entry in the forward walk matrix \mathbf{W} . Applying identical algebraic rearrangement using the column-normalised forward matrix

$$w_{k,j} = \frac{\Delta_{k,j}}{C_j}:$$

$$\begin{aligned} g_j &= \frac{1}{\sum_{k>j} w_{k,j}} + \sum_{k>j} \frac{w_{k,j}}{\sum_{m>j} w_{m,j}} g_k \\ &\leq \frac{1}{\sum_{k>j} w_{k,j}} + \max_{k>j} g_k \end{aligned}$$

Because node j is not the sink, $\sum_{k>j} \Delta_{k,j} \geq c$. Substituting $W_{fwd}(j \rightarrow k)$:

$$\frac{1}{\sum_{k>j} w_{k,j}} = \frac{C_j}{\sum_{k>j} \Delta_{k,j}} \leq \frac{C_j}{c}$$

which implies that

$$g_j \leq \frac{C_j}{c} + \max_{k>j} g_k$$

Recursively expanding this yields $\max_j g_j = g_1 \leq \frac{1}{c} \sum_{k=1}^{n-1} C_k$. Because Δ is row-stochastic, the sum of all column normalisation constants C_k cannot exceed the total number of rows n . Thus, $g_1 \leq \frac{n}{c}$. Similar to our analysis above, by Markov's Inequality, $T_{fwd} \leq \frac{4n}{c}$. \square

Hence, observe that for any $n \geq 2$, instantaneous mixing is impossible. Therefore, both $T_{bwd} \geq 1$ and $T_{fwd} \geq 1$. Utilising Lemma C.1, we can bound the backward time linearly by the forward time:

$$T_{bwd} \leq \frac{4n}{c} \cdot 1 \leq \frac{4n}{c} \cdot T_{fwd}$$

We now apply this relation to our General Fidelity-Mixing Bound (Corollary 6.1), substituting $f(T_{fwd}) = \frac{4n}{c} \cdot T_{fwd}$:

$$F(G) \leq \frac{8}{3} \left(\frac{4n}{c} \cdot T_{fwd}(G) \right) = \frac{32n}{3c} \cdot T_{fwd}(G)$$

and we are done. \square