

Du processus de diffusion sur les graphes aux GNN's ou comment dépasser le Message Passing

Séminaire TALia du 1 avril 2022

[Source](#)

arXiv > cs > arXiv:2106.10934

Computer Science > Machine Learning

[Submitted on 21 Jun 2021 (v1), last revised 22 Sep 2021 (this version, v2)]

GRAND: Graph Neural Diffusion


Benjamin Paul Chamberlain, James Rowbottom, Maria Gorinova, Stefan Webb,
Emanuele Rossi, Michael M. Bronstein

Plan

1. Limitations des GNN
2. L'équation de diffusion dans \mathbb{R}^d
3. L'équation de diffusion sur une variété riemannienne
4. L'équation de diffusion sur un graphe
5. Le modèle Graph Neural Diffusion
6. Comparaison avec les GNN usuels
7. Conclusion

tds






Published in Towards Data Science



Michael Bronstein

Mar 3 · 29 min read ★ · [Listen](#)

Following

A NEW BLUEPRINT FOR LEARNING ON GRAPHS

Graph Neural Networks beyond Weisfeiler-Lehman and vanilla Message Passing

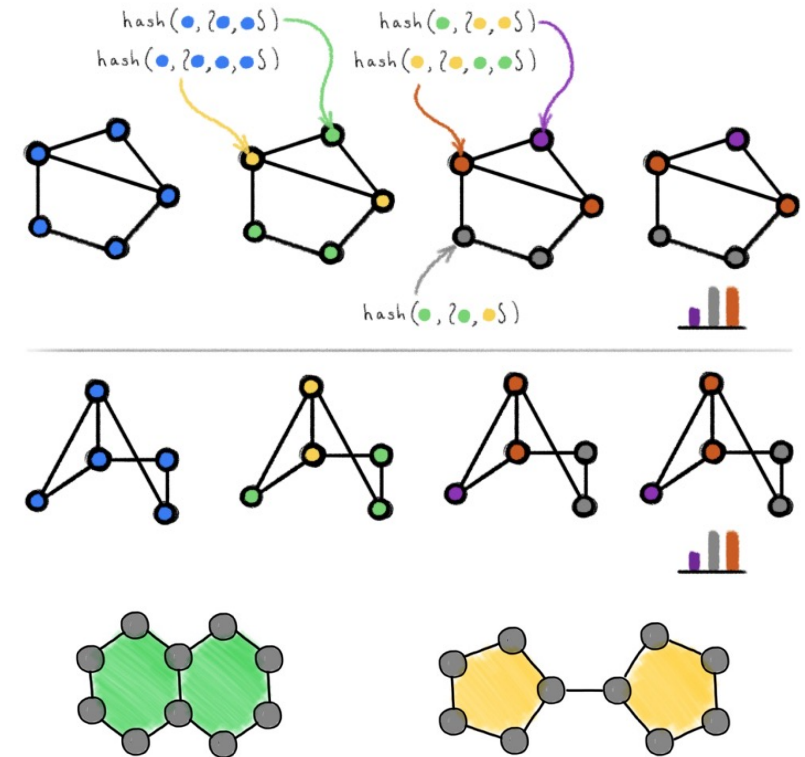
Physics-inspired continuous learning models on graphs allow to overcome the limitations of traditional GNNs

[référence](#)

Message Passing : limitations, difficultés et choix arbitraires

$$\begin{aligned} \mathbf{h}_u^{(k+1)} &= \text{UPDATE}^{(k)} \left(\mathbf{h}_u^{(k)}, \text{AGGREGATE}^{(k)}(\{\mathbf{h}_v^{(k)}, \forall v \in \mathcal{N}(u)\}) \right) \\ &= \text{UPDATE}^{(k)} \left(\mathbf{h}_u^{(k)}, \mathbf{m}_{\mathcal{N}(u)}^{(k)} \right), \end{aligned}$$

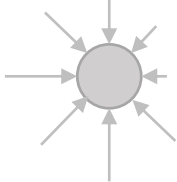
- Le paradigme **message-passing** des GNN (approches spatiales, spectrales) possède toutes les limitations de l'algorithme WL.
- **Over-smoothing** : les GNN trop profonds ne parviennent plus à distinguer les nœuds individuellement (filtre passe-bas).
- **Over-squashing** : lorsque l'information de nœuds lointains doit passer par des goulets d'étranglement
- **Hétérophile** difficile à prendre en compte (filtres passe-haut)
- L'invariance sous la renumérotation des nœuds implique l'**isotropie** dans la propagation de l'information.
- Inadéquation d'un **espace euclidien** pour représenter la notion de proximité entre nœuds dans un graphe
- Arbitraire des procédures de **rewiring** utilisées pour découpler le graphe de calcul du graphe des données ([référence](#))
- Arbitraire des procédures de **décoration/coloriage** des nœuds ([référence](#))



EDP : un réservoir de concepts, de résultats et de méthodes

Idée générale : envisager un GNN comme une **discrétisation**, dans le temps (couches) et dans l'espace (nœuds), d'un processus de diffusion de l'information (features x_u) défini sur une variété riemannienne.

Équation de **continuité**

$$\frac{\partial x(u, t)}{\partial t} = -\text{div} \mathbf{h}(u, t)$$


$$\nabla \cdot \mathbf{h} = \frac{\partial h_1}{\partial u_1} + \frac{\partial h_2}{\partial u_2} + \frac{\partial h_3}{\partial u_3}$$

Flux (de matière, de chaleur...) lié au gradient

$$\mathbf{h}(u, t) = -g \nabla x(u, t)$$

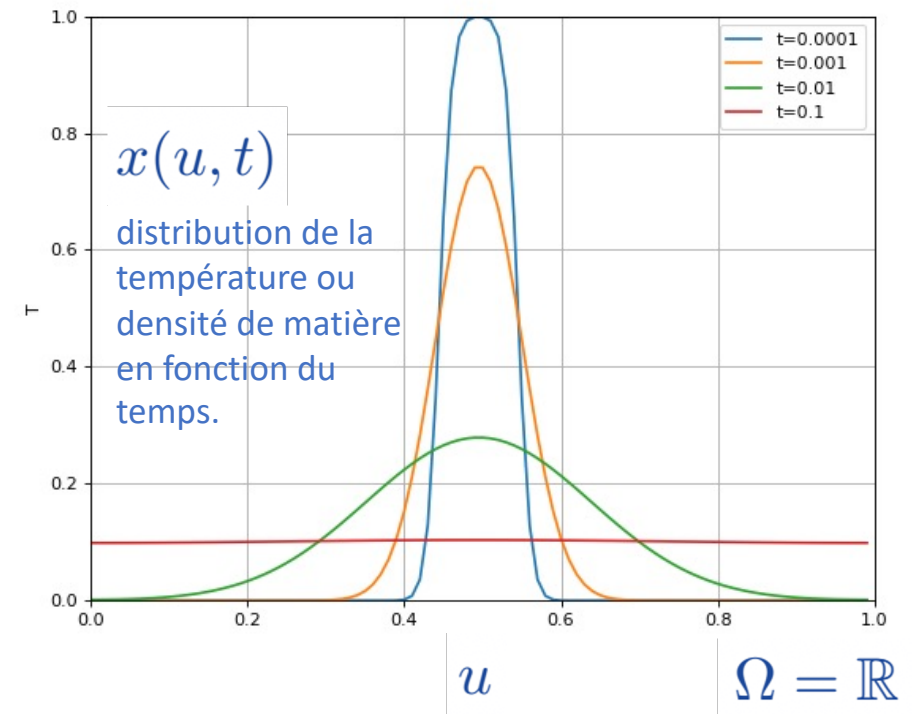
$$g = g(u, x(u, t), t)$$

coefficient de diffusion

Equation de la **diffusion**

$$\frac{\partial x(u, t)}{\partial t} = \text{div} [g(u, x(u, t), t) \nabla x(u, t)]$$

$$x(u, 0) = x_0(t)$$



Solution de l'équation de diffusion lorsque g est constante

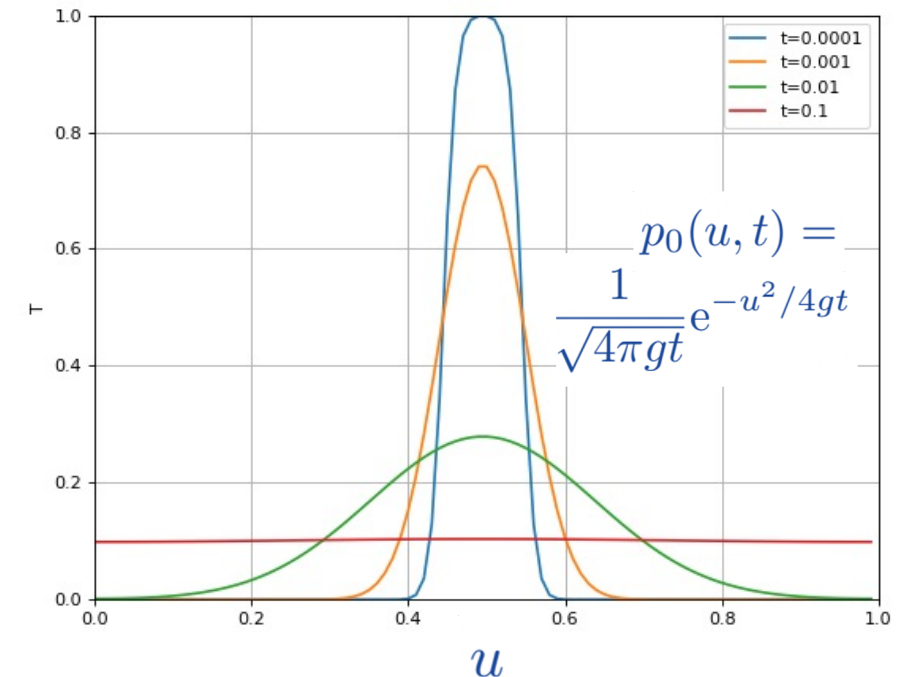
Pour une constante de diffusion g fixe $\frac{\partial x(u, t)}{\partial t} = g \operatorname{div} \nabla x(u, t) = g \Delta x(u, t)$

$$p_0(u, t) = \frac{1}{\sqrt{4\pi gt}} e^{-u^2/4gt}$$

$$\mathbb{E}[u^2] = \int u^2 p_0(u, t) du = 2dgt$$

Diffusion : distance moyenne parcourue en un temps $t \approx \sqrt{t}$

La solution générale s'obtient par linéarité $x(u, t) = \int x_0(v) p_0(u - v, t) dv$



Diffusion sur une variété riemannienne

Graphe = discrétisation spatiale d'une **variété riemannienne** Ω (= surface généralisée avec une notion de distance locale). Permet de rendre naturelles les définitions du gradient ∇ et de la divergence $\nabla \cdot$ sur un graphe.

Fonctions scalaires sur l'espace de base Ω : $x \in \chi(\Omega)$ température, densité, pression, feature,...

Fonctions vectorielles sur l'espace de base Ω : $\mathcal{X} \in \chi(T\Omega)$ flux, direction du vent, ...

Opérateur **gradient** : $\nabla : \chi(\Omega) \rightarrow \chi(T\Omega)$

Opérateur **divergence** : $\text{div} : \chi(T\Omega) \rightarrow \chi(\Omega)$

Pour un espace euclidien $\Omega = \mathbb{R}^d$ on a :
$$\int_{\mathbb{R}^d} \nabla x(u) \cdot \mathcal{X}(u) \, du = \int_{\mathbb{R}^d} x(u) \text{div} \mathcal{X}(u) \, du$$

Pour une variété riemannienne, on a aussi :
(à condition d'utiliser les bonnes définitions)

$$\langle\langle \nabla x, \mathcal{X} \rangle\rangle = \langle x, \text{div} \mathcal{X} \rangle$$

Gradient et divergence sur un graphe

Analogues des définitions pour un graphe $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ $|\mathcal{V}| = n$ $|\mathcal{E}| = e$

features de **nœud** : $x : \mathcal{V} \ni i \rightarrow x_i \in \mathbb{R}$ $x \in \mathbb{R}^n$

features de **lien** : $\mathcal{X} : \mathcal{E} \ni (ij) \rightarrow \mathcal{X}_{ij} \in \mathbb{R}$ $\mathcal{X} \in \mathbb{R}^e$

$$\mathcal{X}_{ji} = -\mathcal{X}_{ij}$$

un nombre attribué à chaque direction j définie un lien du graphe à partir d'un nœud i .

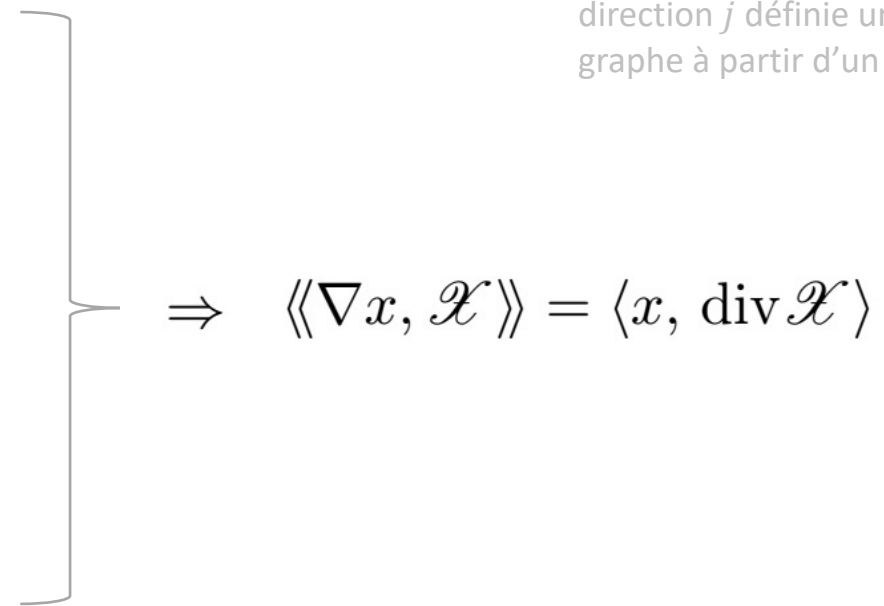
$$(\nabla x)_{ij} := x_j - x_i$$

$$(\text{div}(\mathcal{X}))_i := \sum_{j:(ij) \in \mathcal{E}} \mathcal{X}_{ij} = \sum_{j=1}^n w_{ij} \mathcal{X}_{ij}$$

$$\langle x, y \rangle := \sum_{i \in \mathcal{V}} x_i y_i$$

$$\langle\langle \mathcal{X}, \mathcal{Y} \rangle\rangle := \sum_{(ij) \in \mathcal{E}} \mathcal{X}_{ij} \mathcal{Y}_{ij} = \sum_{i>j} w_{ij} \mathcal{X}_{ij} \mathcal{Y}_{ij}$$

matrice d'adjacence



$$\Rightarrow \langle\langle \nabla x, \mathcal{X} \rangle\rangle = \langle x, \text{div} \mathcal{X} \rangle$$

Diffusion sur un graphe

$$\frac{\partial x(t)}{\partial t} = \text{div} [\mathbf{G}(x(t), t) \nabla x(t)] \quad \text{avec} \quad \mathbf{G}(x, t) = \mathbf{G}(x) := \text{diag}([a(x_i, x_j)]_{(ij) \in \mathcal{E}}) \in \mathbb{R}^{e \times e}$$

En introduisant les définitions de div et ∇ et en définissant une **matrice d'attention** on obtient la forme matricielle de l'équation de diffusion sur un graphe :

$$\mathbf{A}(x) := [\underbrace{a(x_i, x_j)}_{:=0 \text{ si } (ij) \notin \mathcal{E}}]_{i,j \in \mathcal{V}} \in \mathbb{R}^{n \times n}$$

on choisit une matrice stochastique (à droite) : somme par ligne des éléments = 1.

$$\frac{\partial x(t)}{\partial t} = [\mathbf{A}(x(t)) - \mathbf{I}] x(t)$$

Pour une matrice constante $\mathbf{A}(x(t)) = \mathbf{A}$ la solution s'écrit $x(t) = e^{(\mathbf{A} - \mathbf{I})t} x(0)$

Stabilité en ML et dans l'équation de diffusion

Notion de **stabilité en ML** : un algorithme est stable si la fonction de coût ne varie pas beaucoup lorsqu'on remplace une observation par une autre dans le train set.

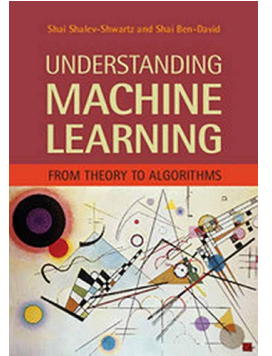
Un **algorithme stable généralise** : l'erreur empirique (sur le train set) est proche de l'erreur théorique (espérance sur la population).

Notion générale de **stabilité pour une EDP**

$$\forall \epsilon > 0, \exists \delta > 0 \quad \text{tel que} \quad |x(0) - \hat{x}(0)| \leq \delta \quad \Rightarrow \quad |x(t) - \hat{x}(t)| \leq \epsilon, \quad \forall t \geq 0$$

Pour la **diffusion sur un graphe** on montre seulement qu'en général

$$\min_{j \in \mathcal{V}} x_j(0) \leq x_i(t) \leq \max_{j \in \mathcal{V}} x_j(0) \quad \forall t \geq 0$$



voir le
théorème 13.2
p139

[PDF](#)

Schémas de résolution explicites et implicites

$$\frac{\partial x(t)}{\partial t} = [\mathbf{A}(x(t)) - \mathbf{I}] x(t)$$

à discrétiser dans le temps

Discrétisation **explicite**

compromis à trouver entre le nombre d'itérations K et le pas τ

$$\frac{x^{(k+1)} - x^{(k)}}{\tau} = \underbrace{[\mathbf{A}(x^{(k)}) - \mathbf{I}]}_{:= \bar{\mathbf{A}}(x^{(k)})} x^{(k)}$$



correspond approximativement à la règle de propagation d'un GAT avec une skip-connexion

$$x^{(k+1)} = [\mathbf{I} + \tau \bar{\mathbf{A}}(x^{(k)})] x^{(k)}$$

schéma de calcul **stable seulement** pour $0 < \tau < 1$

temps discrétisé correspond au n° d'une couche de convolution d'un GNN.

trouver un compromis entre la nb K d'itérations et le pas d'intégration τ .

Discrétisation **implicite**

$$\frac{x^{(k+1)} - x^{(k)}}{\tau} = [\mathbf{A}(x^{(k)}) - \mathbf{I}] x^{(k+1)}$$



$$x^{(k)} = [\mathbf{I} - \tau \bar{\mathbf{A}}(x^{(k)})] x^{(k+1)}$$

schéma de calcul **inconditionnellement stable** pour tout $\tau > 0$

matrice à inverser = prix à payer pour la stabilité

Schémas de résolution multi-étapes

Schéma **Runge-Kutta 4** : utilisation de temps intermédiaire
 Revient à définir des skip-connexions particulières dont l'efficacité a été démontrée en pratique.

$$f(\mathbf{x}, t) = \bar{A}(\mathbf{x}_t) \mathbf{x}_t$$

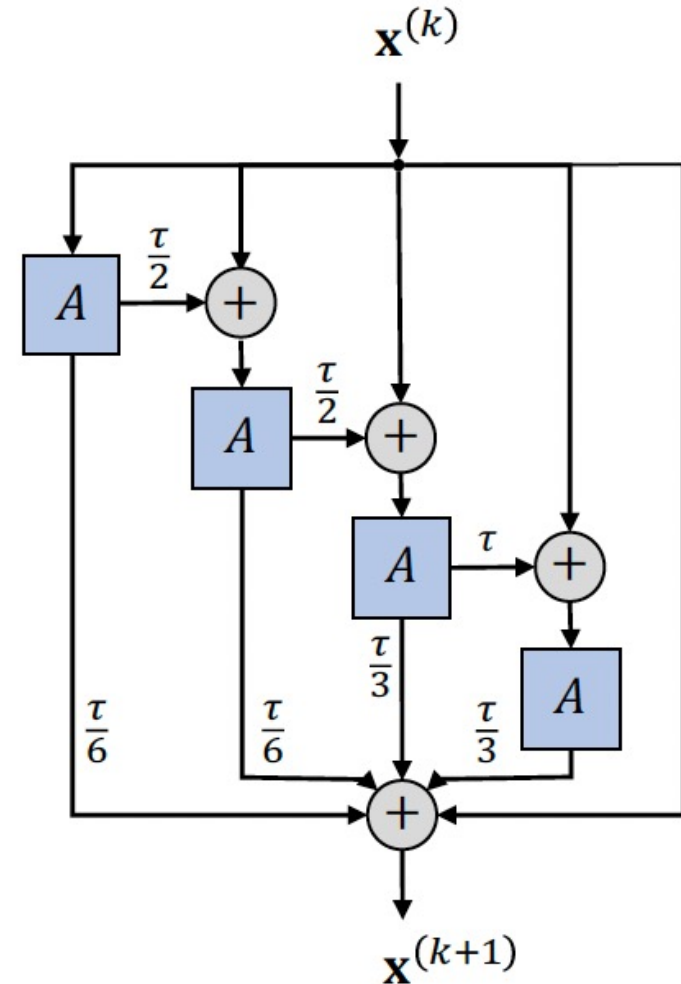
$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \frac{1}{6} \tau (\mathbf{k}_1 + 2\mathbf{k}_2 + 2\mathbf{k}_3 + \mathbf{k}_4)$$

$$\mathbf{k}_1 = f(\mathbf{x}_t, t)$$

$$\mathbf{k}_2 = f(\mathbf{x}_t + \tau \mathbf{k}_1 / 2, t + \tau / 2)$$

$$\mathbf{k}_3 = f(\mathbf{x}_t + \tau \mathbf{k}_2 / 2, t + \tau / 2)$$

$$\mathbf{k}_4 = f(\mathbf{x}_t + \tau \mathbf{k}_3, t + \tau)$$



Le modèle : Graph Neural Diffusion (GRAND)

Comme pour PPNP on va découpler l'apprentissage de représentations locales par un modèle **non-linéaire de type MLP** de leur **propagation** à travers le graphe.

$$Z_{GCN} = \text{softmax} (P\sigma \dots (P\sigma(PXW^{(0)}) W^{(1)}) \dots W^{(L)})$$

$$Z_{PPNP} = \text{softmax} (P_{ppr} f_{\theta}(X)) \quad \text{où} \quad P_{ppr} := \alpha [\mathbb{I} - (1 - \alpha)P]^{-1}$$

Le modèle **GRAND**

- 1 $\mathbf{X}(0) = \phi(\mathbf{X}_{in}) \in \mathbb{R}^{n \times d}$
- 2 $\mathbf{X}(T) = \mathbf{X}(0) + \int_0^T \frac{\partial \mathbf{X}(t)}{\partial t} dt \quad \text{où} \quad \frac{\partial \mathbf{X}(t)}{\partial t} = [\mathbf{A}(\mathbf{X}) - \mathbf{I}] \mathbf{X}$
- 3 $\mathbf{Y} = \psi(\mathbf{X}(T))$

Equation non linéaire si \mathbf{A} dépend de \mathbf{X} .

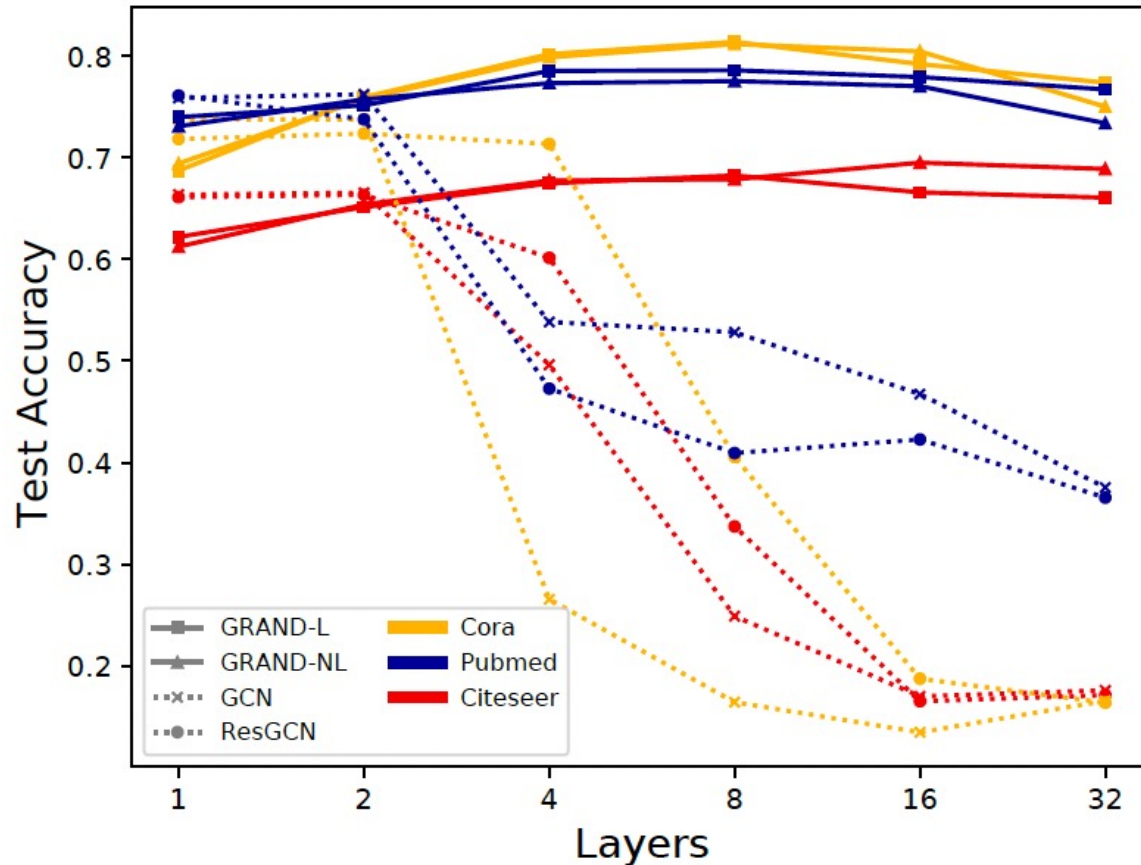
$$\mathbf{A}(\mathbf{X}) := [a(\mathbf{X}_i, \mathbf{X}_j)]_{i,j \in \mathcal{V}}$$

$$a(\mathbf{X}_i, \mathbf{X}_j) := \text{softmax} \left(\frac{(\mathbf{W}_K \mathbf{X}_i)^\top \mathbf{W}_Q \mathbf{X}_j}{d_k} \right)$$

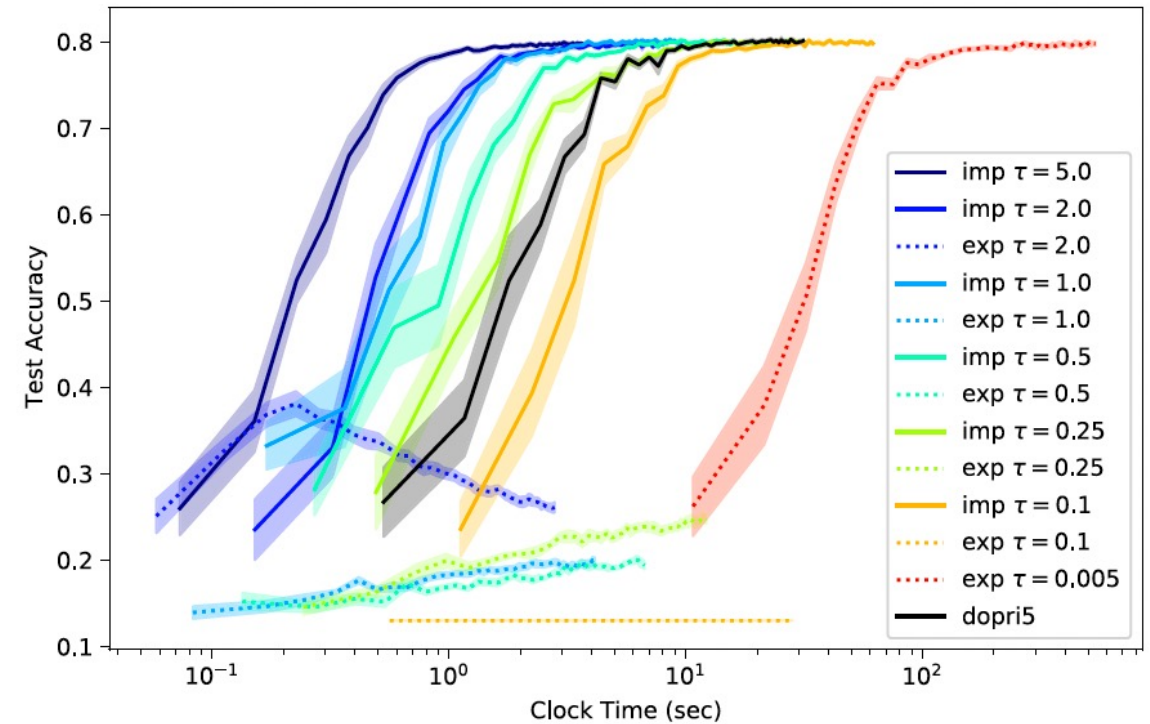
Choix à effectuer pour la propagation par diffusion :

- Matrice de diffusion $\mathbf{G}(\mathbf{X})$ \leftrightarrow matrice d'attention $\mathbf{A}(\mathbf{X})$.
- Choix d'une méthode de discrétisation temporelle de l'équation de diffusion.

Résultats et comparatifs avec les GNN classiques



Robustesse à l'oversmoothing en fonction du nombre K de couches/pas d'intégration avec $\tau=1$.



Stabilité des schémas implicites et explicites en fonction de la taille du pas d'intégration τ

Conclusion

- **GRAND** est une incarnation concrète de la stratégie qui consiste à transférer des techniques, mathématiques et numériques, développée pour la physique à l'apprentissage de représentations pour les graphes.
- C'est une partie d'une ambition beaucoup plus vaste qui entend convoquer la théorie des EDP, le control optimal, la géométrie différentielle, la topologie algébrique pour refonder les GNN's.
- **GRAND** propose une méthode systématique pour construire de nouveaux modèles d'apprentissage de représentations pour les graphes, avec des méthodes adaptatives par exemple.
- **GRAND** permet d'interpréter certains modèles existants comme des cas particuliers.
- Le choix d'une méthode d'intégration numérique (discrétisation) revient à définir des connexions résiduelles (skip-connexions) dont on a de bonnes raisons de penser qu'elles sont utiles.
- Les résultats de stabilité permettent de construire des modèles très profonds.
- Le modèle **GRAND** n'utilise qu'un nombre relativement faible de paramètres, ceux qui définissent l'unique mécanisme d'attention (au lieu d'un mécanisme par couche).

Question à 1 millions d'€ : S'agit-il d'un simple « rebranding chic » d'idées anciennes ? Une montagne mathématique a-t-elle accouché d'une souris au niveau des applications ?