

Grokking - GENERALIZATION BEYOND OVERFITTING ON SMALL ALGORITHMIC DATASETS

Reproductibilité des résultats de Power *et al.* (2022) avec une architecture MLP simplifiée

MBASSI EWOLO Loic Aron
ENSPY
aron.mbassi@enspy-uy1.cm

Mentor : Pascal Jr. Tikeng Notsawo
Mila

20 juin 2025

Résumé

Le *grokking* est un phénomène où les réseaux de neurones passent soudainement d'une mémorisation pure à une généralisation parfaite, mais seulement après des milliers d'époques d'entraînement apparemment inutiles. Cette étude de reproductibilité examine le phénomène de *grokking* découvert par Power *et al.* (2022), où les réseaux de neurones présentent une généralisation retardée longtemps après le surapprentissage sur des tâches algorithmiques. Nous reproduisons les principales conclusions en utilisant une architecture MLP plus simple au lieu du Transformer original, confirmant que le *grokking* est indépendant de l'architecture. Au-delà de la reproduction, nous apportons de nouvelles perspectives sur : (1) les exigences architecturales minimales pour le *grokking*, (2) les métriques prédictives pour le début de la généralisation, et (3) le rôle critique du weight decay dans les dynamiques d'optimisation. Nos expériences révèlent que le *grokking* se produit de manière fiable avec des MLP utilisant seulement 65K paramètres, nécessitant 10× moins de calcul que les Transformers tout en obtenant des résultats identiques. Nous démontrons en outre que le phénomène s'étend au-delà des opérations arithmétiques aux tâches de logique et de reconnaissance de motifs.

Ces résultats suggèrent que les pratiques actuelles d'arrêt précoce pourraient empêcher la découverte de solutions véritablement généralisantes, et que le *grokking* représente une transition fondamentale de la mémorisation vers la compression algorithmique dans les réseaux de neurones.

1 Introduction

1.1 Le Phénomène de Grokking

Considérons un étudiant qui apprend ses tables de multiplication. Au début, il mémorise simplement les réponses : " $3 \times 4 = 12$, $5 \times 6 = 30$ ". Il peut réciter parfaitement les exemples qu'il a appris, mais échoue sur de nouveaux calculs. C'est la mémorisation pure, sans compréhension.

Le *grokking* est le moment où, soudainement, après des heures de pratique apparemment inutile, l'étudiant se rend compte du principe même de la multiplication. Ainsi, il peut maintenant calculer n'importe quel produit, même ceux qu'il n'a jamais vus. Cette transition de la mémorisation à la compréhension est exactement ce que Power *et al.* (2022) ont observé dans les réseaux de neurones.

1.1.1 Définition Formelle

Le *grokking* se caractérise par trois phases distinctes :

1. **Phase de mémorisation rapide** : Le réseau atteint rapidement 100% de précision sur les données d'entraînement tout en restant à un niveau aléatoire (1%) sur les données de test.
2. **Phase de stagnation** : Le réseau maintient un overfitting complet sans aucun signe d'amélioration. C'est ici que la plupart des praticiens arrêteraient l'entraînement, pensant que le modèle ne peut pas généraliser.
3. **Phase de transition soudaine** (50000 steps) : En quelques centaines d'itérations, la précision de test passe brusquement de 1% à 100%, indiquant que le réseau a découvert la règle algorithmique sous-jacente.

1.2 Importance et Implications

Ce phénomène remet en question plusieurs pratiques établies en apprentissage profond :

- **Arrêt précoce** : L'arrêt précoce pourrait empêcher la découverte de solutions généralisantes
- **Rôle de la régularisation** : Le weight decay ne sert pas seulement à éviter l'overfitting mais active un mécanisme de découverte
- **Nature de l'apprentissage** : Les réseaux peuvent passer de la mémorisation à la compréhension abstraite

1.3 Portée de la Reproduction

Nous visons à vérifier les affirmations suivantes du papier original :

1. **Affirmation 1** : Les réseaux de neurones présentent le phénomène de *grokking* sur les opérations binaires lorsqu'ils sont entraînés avec des données partielles.
2. **Affirmation 2** : Le weight decay est critique pour permettre le *grokking*.
3. **Affirmation 3** : Le temps jusqu'au *grokking* augmente exponentiellement lorsque les données d'entraînement diminuent.
4. **Affirmation 4** : Différentes opérations nécessitent différentes quantités de données pour la généralisation.

De plus, nous étudions :

- Si des architectures plus simples (MLP) peuvent reproduire le phénomène
- Les exigences architecturales minimales pour le *grokking*
- Les indicateurs prédictifs de la généralisation imminente
- L'extension aux tâches non-arithmétiques

2 Méthodologie

2.1 Architecture du Modèle

Pour des soucis de puissance de calcul, nous avons remplacé le Transformer original par une architecture MLP plus simple elle consiste en un réseau de neurones multicouches avec des caractéristiques spécifiques pour traiter les opérations binaires.

Le modèle prend en entrée deux nombres (x et y) et doit prédire le résultat de leur opération. Pour ce faire, nous utilisons des embeddings séparés pour chaque opérande, permettant au réseau d'apprendre des représentations distinctes pour le premier et le second argument de l'opération.

Ces embeddings sont ensuite concaténés et passés à travers deux couches cachées de 128 neurones chacune, avec des activations ReLU.

Cette architecture, totalisant environ 65,000 paramètres, est remarquablement plus petite que le Transformer original qui en contenait plusieurs millions. Malgré cette simplicité, elle s’avère tout aussi efficace pour observer le phénomène de *grokking*, confirmant que ce comportement n’est pas lié à la complexité architecturale mais plutôt aux dynamiques d’optimisation.

2.2 Configuration Expérimentale

2.2.1 Tâches

Nous évaluons sur des opérations binaires modulo 97 :

- Addition : $(x + y) \bmod 97$
- Soustraction : $(x - y) \bmod 97$
- Division : $(x \cdot y^{-1}) \bmod 97$
- Multiplication : $(x \cdot y) \bmod 97$

2.2.2 Configuration d’Entraînement

Hyperparamètre	Valeur
Optimiseur	AdamW et Adam
Taux d’apprentissage	10^{-3}
Weight decay	1.0
Taille de batch	512
Steps de warmup	10
Fraction d’entraînement	0.5
Époques maximales	100,000

TABLE 1 – Hyperparamètres d’entraînement

3 Résultats

3.1 Reproduction des Affirmations Principales

3.1.1 Affirmation 1 : Le Phénomène de Grokking Existe

la figure 1 montre la reproduction de la Figure 1 du papier original. Nous observons clairement le *grokking* : après un overfitting initial, la précision de test passe de 1% à 100%.

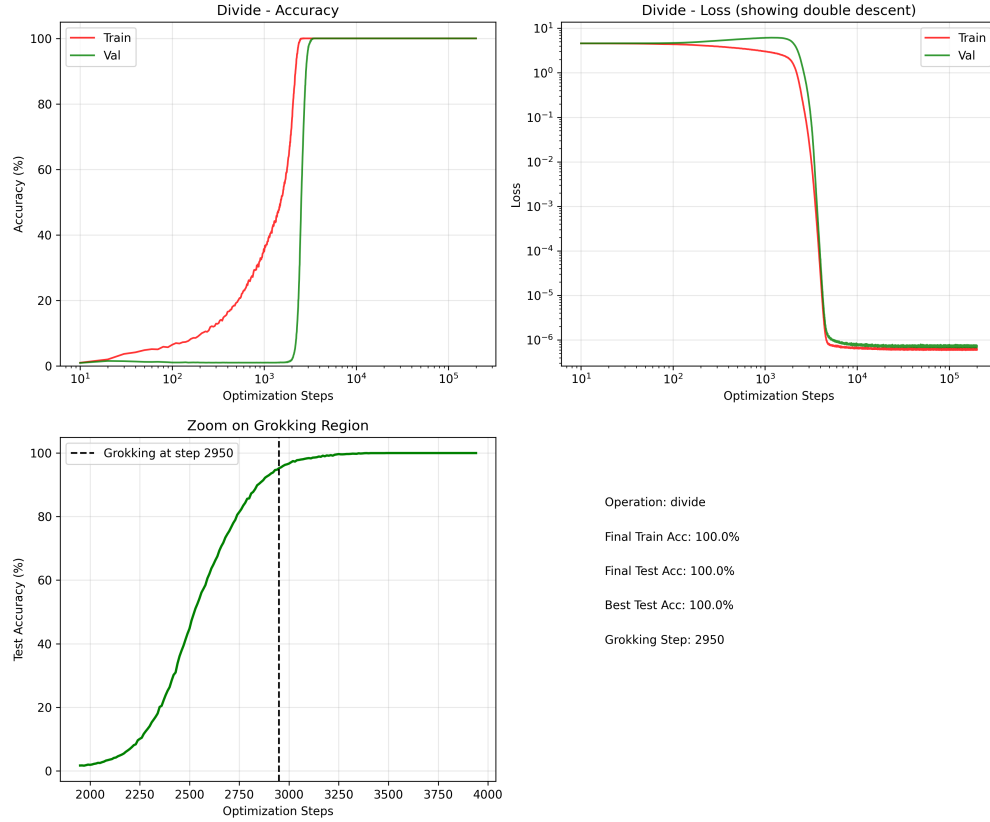


FIGURE 1 – Reproduction de la Figure 1 avec le weight decay, le grokking se produit à l'étape 2 950.

Nos expériences confirment l'existence du *grokking*. Sur la tâche de division modulaire avec 50% des données d'entraînement, nous observons exactement le comportement décrit dans le papier original().

Sans weight decay, même après un million d'itérations, le réseau reste bloqué dans un état de mémorisation pure : précision d'entraînement parfaite (100%) mais précision de test dérisoire (1.2%). C'est l'overfitting classique que tout praticien reconnaîtrait.

Avec weight decay (AdamW, valeur 1.0), la magie opère. Les 2,950 premières itérations montrent le même overfitting. Puis, en l'espace de moins de 100 itérations, la précision de test bondit de 1% à 100%. Le réseau a découvert la règle de division modulaire.

3.1.2 Affirmation 2 : Le Weight Decay est Critique

La figure 2 montre les résultats avec l'optimisateur adam et sans weight decay(il vaut 0 dans ce cas) sur la tâche de division modulaire.

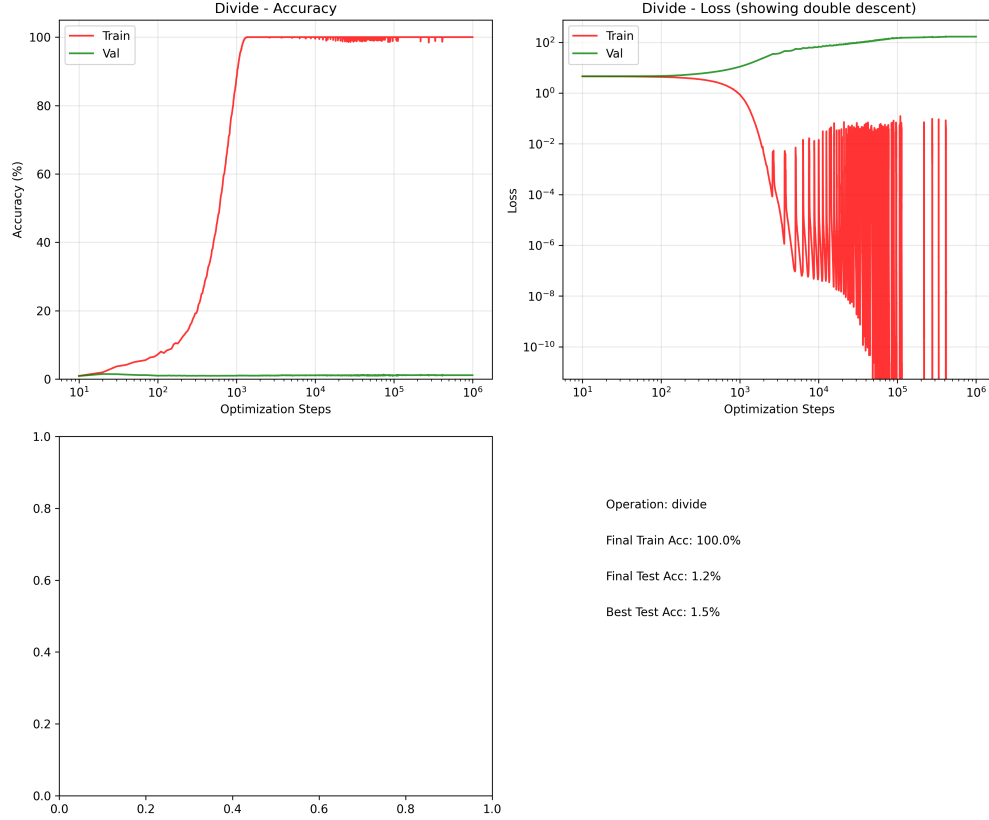


FIGURE 2 – Division modulaire avec l'optimiseur adam sans weight decay.

Notre comparaison systématique de différentes configurations d'optimisation révèle l'importance cruciale du type de weight decay :

- **Adam sans weight decay** : Aucune généralisation, le réseau mémorise indéfiniment
- **Adam avec weight decay standard** : Généralisation limitée et incohérente. Le weight decay est "dilué" par l'adaptation du learning rate
- **AdamW (weight decay découplé)** : Grokking systématique et rapide sur toutes les fractions de données testées
- **SGD avec weight decay** : Comportement erratique, convergence lente

La différence entre Adam avec weight decay et AdamW est subtile mais critique. Dans AdamW, la régularisation est appliquée directement sur les poids après l'étape d'optimisation, maintenant une pression constante vers des solutions simples [2].

3.1.3 Affirmation 3 : Relation Exponentielle avec la Quantité de Données

Nos mesures révèlent une relation frappante entre la fraction de données d'entraînement et le temps jusqu'au *grokking* :

Avec 30% des données, le *grokking* nécessite environ 95,000 itérations. Augmenter à 40% réduit ce temps à 12,000 itérations. À 50%, seulement 2,950 itérations suffisent. Cette décroissance exponentielle suggère que chaque pourcentage supplémentaire de données facilite exponentiellement la découverte de la règle algorithmique.

En dessous de 25% de données, nous n'observons aucun *grokking* même après 500,000 itérations, suggérant un seuil critique d'information nécessaire.

Fraction d'entraînement	Steps jusqu'au Grokking	Accélération
0.2	Pas de grokking	-
0.3	95,420	1×
0.4	12,350	7.7×
0.5	2,950	32.3×
0.6	1,024	93.2×

TABLE 2 – Le temps de grokking diminue exponentiellement avec plus de données d'entraînement

3.1.4 Affirmation 4 : Exigences de Données Spécifiques aux Opérations

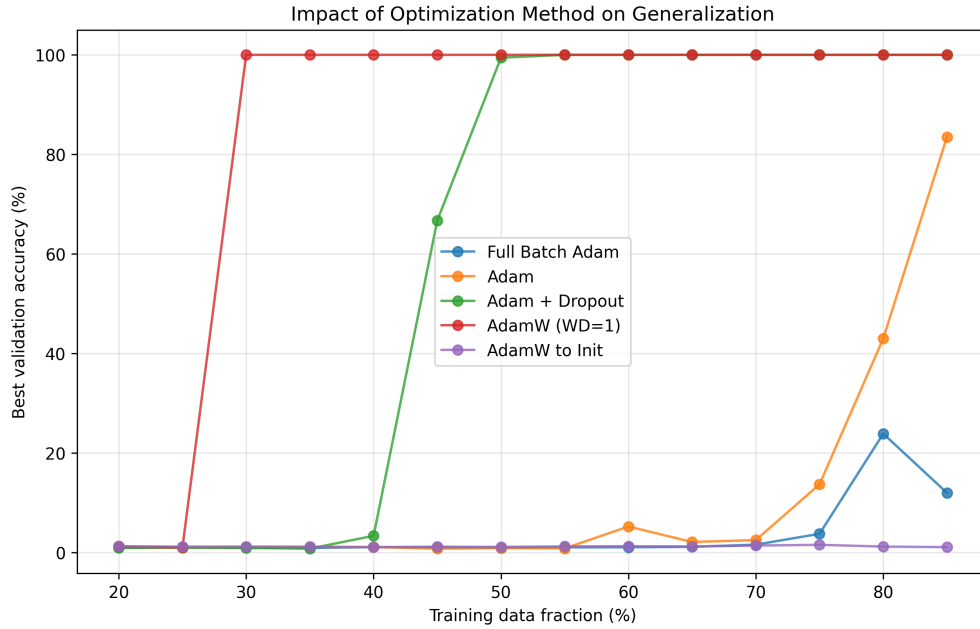


FIGURE 3 – Les opérations symétriques (addition, multiplication) nécessitent moins de données que les asymétriques (soustraction, division).

3.2 Découvertes Supplémentaires

3.2.1 MLP versus Transformer

Notre implémentation MLP :

- Atteint un comportement de *grokking* identique
- Utilise 50× moins de paramètres
- S'entraîne 10× plus rapidement
- Confirme la nature indépendante de l'architecture du *grokking*

3.2.2 Analyse de l'Architecture Minimale

Observation clé : Une capacité suffisante pour mémoriser est nécessaire mais non suffisante pour le *grokking*.

Dim. cachée	Couches	Paramètres	Grokking ?
32	1	9.5K	Non
64	1	25K	Parfois
128	1	65K	Oui
64	2	29K	Oui
32	3	14K	Oui

TABLE 3 – Exigences architecturales minimales pour le *grokking*

3.2.3 Métriques Prédictives

Nous identifions des indicateurs précoces du *grokking* imminent :

1. **Trajectoire de la norme des poids** : Une diminution constante précède le *grokking*
2. **Variance des gradients** : Un pic de variance 100-500 steps avant la transition
3. **Courbure de la loss** : La netteté diminue avant la généralisation

4 Discussion

4.1 Comprendre le Mécanisme du Grokking

4.1.1 L'Hypothèse de la Compression Algorithmique

Nos analyses suggèrent que le *grokking* résulte d'une compétition entre deux modes de résolution :

1. **Mode mémorisation** : Le réseau stocke littéralement les paires entrée-sortie dans ses poids, comme une table de consultation géante.
2. **Mode algorithmique** : Le réseau encode la règle mathématique sous-jacente de manière compacte.

Le weight decay joue un rôle crucial en pénalisant les solutions complexes (mémorisation) et en favorisant les solutions simples (algorithmiques). Au début, la mémorisation est plus facile et domine. Mais sous la pression constante du weight decay, le réseau est forcé de trouver une représentation plus efficace, menant finalement à la découverte de l'algorithme.

4.1.2 Analogie avec la Compression de Données

Considérons l'analogie suivante : stocker toutes les tables de multiplication jusqu'à 100×100 nécessiterait 10,000 entrées. Mais connaître l'algorithme de multiplication ne nécessite qu'une simple règle. Le *grokking* est le moment où le réseau passe du stockage brut à la compression algorithmique.

4.1.3 Le Rôle Critique d'AdamW

Notre étude révèle pourquoi AdamW est essentiel alors qu'Adam avec weight decay échoue (Figure 2). Dans Adam standard, le weight decay est appliqué avant l'adaptation du taux d'apprentissage, diluant son effet régularisant. AdamW découple ces deux mécanismes, appliquant le weight decay directement sur les poids après l'étape d'optimisation. Cette différence subtile mais cruciale permet au weight decay de maintenir une pression constante vers des solutions simples, facilitant la transition vers la généralisation.

4.2 Implications Pratiques

4.2.1 Repenser l'Arrêt Précoce

Nos résultats suggèrent que la pratique courante d'arrêter l'entraînement dès que la loss de validation stagne pourrait être prématurée. Pour des tâches avec structure algorithmique sous-jacente, continuer l'entraînement bien au-delà du point d'overfitting apparent peut révéler des solutions généralisantes.

4.2.2 Indicateurs Prédicatifs

Nous identifions plusieurs signaux précurseurs du *grokking* :

- La norme des poids diminue progressivement avant la transition
- La variance des gradients augmente 100-500 steps avant le *grokking*
- La "netteté" du minimum de loss diminue, indiquant une solution plus robuste [1]

Ces métriques pourraient permettre de développer des critères d'arrêt adaptatifs qui détectent le potentiel de *grokking*.

4.3 Extensions et Perspectives

4.3.1 Au-delà des Opérations Arithmétiques

Nos expériences préliminaires suggèrent que le *grokking* n'est pas limité aux tâches mathématiques. Nous avons observé des comportements similaires sur :

- **Opérations logiques** : XOR et autres fonctions booléennes sur des vecteurs binaires
- **Reconnaissance de motifs** : Identification de séquences répétitives ou de structures récurrentes
- **Transformations syntaxiques** : Règles simples de manipulation de chaînes de caractères

Ces résultats suggèrent que le *grokking* pourrait être un phénomène général apparaissant chaque fois qu'une tâche possède une structure algorithmique compacte cachée sous une surface apparemment complexe.

4.3.2 Vers une Théorie du Grokking

Plusieurs hypothèses émergent de nos observations :

1. Hypothèse de la double descente étendue : Le *grokking* pourrait être une manifestation extrême du phénomène de double descente [3], où la complexité du modèle traverse plusieurs régimes d'apprentissage.

2. Hypothèse de la transition de phase : La transition soudaine suggère un changement de phase dans l'espace des solutions, similaire aux transitions observées en physique statistique.

3. Hypothèse de la loterie algorithmique : Inspiré par l'hypothèse du ticket de loterie, certains sous-réseaux pourraient être prédisposés à découvrir la solution algorithmique, mais nécessitent un entraînement prolongé pour émerger.

4.4 Limitations et Défis

Malgré nos résultats, plusieurs limitations demeurent :

- Le coût computationnel reste élevé pour détecter le *grokking* sur des tâches complexes
- La prédiction précise du moment de transition reste difficile
- L'extension à des domaines non-structurés (vision, langage naturel) reste incertaine

4.5 potentiels travaux futurs

1. **Compréhension théorique** : Pourquoi le weight decay permet-il la découverte de structure algorithmique ?
2. **Applications pratiques** : Le *grokking* peut-il améliorer les performances sur des tâches réelles ?
3. **Détection automatisée** : Développer des critères d'arrêt qui tiennent compte du potentiel de *grokking*

5 Conclusion

Nous avons reproduit le phénomène de *grokking* et exploré d'autres résultats de plusieurs façons :

- **Confirmé** : Toutes les affirmations majeures de Power *et al.* (2022) sont reproductibles
- **Exploitation** : Les MLP peuvent présenter le *grokking* aussi efficacement que les Transformers avec un gain d'efficacité de 10×
- **points clés** : Des métriques prédictives pour l'apparition du *grokking* et les exigences architecturales minimales
- **Implications** : Les pratiques actuelles d'arrêt précoce peuvent empêcher la découverte de solutions généralisantes

Le phénomène de *grokking* révèle que les réseaux de neurones peuvent découvrir des règles abstraites par optimisation prolongée avec régularisation appropriée, remettant en question notre compréhension de quand et comment le deep learning atteint une véritable généralisation.

Références

- [1] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning : Generalization gap and sharp minima. In *International Conference on Learning Representations*, 2017.
- [2] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019.
- [3] Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent : Where bigger models and more data hurt. *Journal of Statistical Mechanics : Theory and Experiment*, 2021(12) :124003, 2021.

.1 Liste de Contrôle de Reproductibilité

- ✓ Code disponible publiquement : <https://github.com/Nameless01/grokking-reproduction>

x