

La Puce Analogique que les États-Unis ne Peuvent Arrêter

Catégorie : Intelligence Artificielle | Lecture : 13 min | Publié le : 23/03/2026 | Auteur : Ayi NEDJIMI

Puce analogique RRAM de Pékin : 1 000× les performances GPU Nvidia, 100× moins énergivore. Analyse technique, géopolitique et impact des sanctions.

En octobre 2025, une annonce provenant de l'Université de Pékin a secoué le monde technologique mondial. Des chercheurs chinois, menés par le professeur Sun Zhong de l'Institut d'Intelligence Artificielle, ont dévoilé une puce analogique révolutionnaire basée sur la mémoire résistive à accès aléatoire (RRAM). Cette innovation, publiée dans la prestigieuse revue *Nature Electronics*, promet des performances jusqu'à mille fois supérieures aux meilleurs processeurs graphiques actuels, tout en consommant cent fois moins d'énergie. Cette percée symbolise un changement de paradigme fondamental remettant en question des décennies de domination du calcul numérique. Elle intervient dans un contexte géopolitique tendu où les États-Unis ont imposé des sanctions sévères sur les semi-conducteurs, et paradoxalement, ces mêmes restrictions semblent avoir catalysé une innovation radicale qui pourrait redéfinir l'ensemble de l'industrie des puces pour l'intelligence artificielle.

L'ironie de la situation n'échappe à personne. Alors que l'Occident s'est concentré sur la course aux nanomètres et aux architectures numériques toujours plus complexes, la Chine a emprunté un chemin alternatif, revisitant une technologie que le monde numérique avait abandonnée depuis des décennies : le calcul analogique. Une décision forcée par les contraintes, qui s'avère aujourd'hui être un raccourci vers une percée que personne n'avait anticipée. Cet article analyse en profondeur les dimensions techniques, géopolitiques et économiques de cette innovation, ses applications concrètes, et ce qu'elle signifie pour l'avenir de la compétition mondiale en matière d'**intelligence artificielle**.

Comprendre le Calcul Analogique — Un Retour aux Sources

L'Histoire Oubliée du Calcul Analogique

Pour comprendre l'importance de cette percée, il faut remonter aux origines de l'informatique. Avant l'ère des ordinateurs numériques, le calcul analogique dominait le paysage technologique. Les premiers calculateurs utilisaient des phénomènes physiques continus — courants électriques, mouvements mécaniques — pour effectuer des opérations mathématiques. Dans les années 1940 et 1950, les calculateurs analogiques étaient omniprésents dans les laboratoires de recherche et les installations militaires. Ils excellaient dans la résolution d'équations différentielles et la simulation de systèmes physiques complexes. Cependant, leur talon d'Achille résidait dans leur précision limitée et leur sensibilité au bruit.

L'avènement du transistor et la miniaturisation progressive des circuits numériques ont sonné le glas de cette ère. Les ordinateurs numériques, représentant l'information en séquences discrètes de zéros et de uns, offraient une précision parfaite et une reproductibilité sans faille. Le monde informatique s'est engagé sur la voie du numérique, suivant la loi de Moore qui promettait un doublement régulier de la densité des transistors. La technologie analogique semblait définitivement reléguée aux musées de l'histoire de l'informatique.

Les Limites du Paradigme Numérique

Pendant des décennies, le calcul numérique a régné en maître. Cependant, les signaux d'essoufflement se sont multipliés. La loi de Moore montre des signes évidents de ralentissement — la miniaturisation des transistors approche des limites atomiques, rendant chaque nouvelle génération de puces exponentiellement plus coûteuse. Plus problématique encore, l'architecture de von Neumann souffre d'un goulot d'étranglement structurel : le processeur et la mémoire étant séparés physiquement, les transferts constants de données entre ces deux composants consomment une quantité disproportionnée d'énergie.

L'explosion de l'intelligence artificielle a exacerbé ces problèmes. Former un modèle comme GPT-4 nécessite des milliers de GPU fonctionnant pendant des semaines, consommant autant d'électricité qu'une petite ville. Jensen Huang, PDG de Nvidia, a lui-même reconnu que l'électricité — plus que le silicium — pourrait devenir le facteur limitant dans la course à l'IA. C'est précisément cette convergence de crises qui a rendu le moment propice à une renaissance du calcul analogique.

Le Calcul Analogique : Une Renaissance Inattendue

Contrairement aux processeurs numériques qui décomposent l'information en bits discrets, les systèmes analogiques traitent l'information sous forme de signaux continus. Cette approche présente des avantages fondamentaux pour les opérations matricielles qui constituent le cœur des algorithmes d'intelligence artificielle. Imaginez deux nageurs dans une piscine : le premier (calcul numérique) sort de l'eau tous les deux mètres pour courir quelques pas avant de replonger — gaspillant énormément d'énergie dans ces transitions ; le second (calcul analogique) glisse naturellement d'un bout à l'autre du bassin, exploitant les propriétés physiques naturelles. Le problème historique de la faible précision analogique semblait insurmontable. C'est précisément ce *problème centenaire* que l'équipe de l'Université de Pékin prétend avoir résolu.

La Percée de l'Université de Pékin — Décryptage Technique

La Mémoire RRAM : Le Cœur de l'Innovation

Au centre de cette percée se trouve la **RRAM (Resistive Random-Access Memory)**. Contrairement aux mémoires traditionnelles qui stockent l'information sous forme de charges électriques, la RRAM utilise des variations de résistance électrique pour encoder les données. Cette caractéristique unique permet de stocker et traiter l'information simultanément, éliminant ainsi le goulot d'étranglement de l'architecture de von Neumann.

La RRAM fonctionne en modifiant la conductivité d'un matériau oxyde métallique situé entre deux électrodes. En appliquant différentes tensions, on peut créer ou détruire des filaments conducteurs dans cet oxyde, modifiant ainsi sa résistance — et cette résistance peut prendre des valeurs continues, contrairement aux mémoires numériques limitées à deux états. L'équipe de Sun Zhong a organisé ces cellules RRAM en réseaux matriciels appelés *crossbar arrays*. Dans cette configuration, les opérations de multiplication matricielle fondamentales pour l'IA s'effectuent directement par les lois physiques : la loi d'Ohm pour la multiplication (courant = tension × conductance) et la loi de Kirchhoff pour l'addition (les courants s'additionnent le long des colonnes). Vous pouvez consulter l'étude originale sur Nature Electronics.

Résoudre le Problème Centenaire de la Précision

L'innovation clé réside dans une **architecture à double circuit**. Le premier circuit effectue des calculs approximatifs rapides, exploitant la vitesse inhérente du calcul analogique. Le second circuit affine ces résultats par itérations successives, corrigeant les erreurs et atteignant progressivement la précision souhaitée. Cette approche hybride combine le meilleur des deux mondes. Les chercheurs affirment avoir atteint une précision en virgule fixe de **24 bits** — comparable aux systèmes numériques et inédite pour un système analogique. L'amélioration par rapport aux systèmes analogiques précédents est stupéfiante : cinq ordres de grandeur, soit une multiplication par cent mille de la précision.

Performances et Implications

Les chiffres sont vertigineux. La puce RRAM a démontré un débit de calcul plus de **1 000 fois supérieur** au Nvidia H100, avec une efficacité énergétique **100 fois supérieure**. Une tâche nécessitant une journée entière sur un GPU moderne pourrait être accomplie en environ une minute, tout en consommant une fraction de l'énergie. Les implications pour les centres de données IA sont considérables : réduction massive des coûts d'exploitation, diminution de l'empreinte carbone, et possibilité de déployer des capacités de calcul dans des environnements précédemment inaccessibles. Aspect crucial pour la commercialisation : la puce a été fabriquée avec des **processus de production commerciaux standard**, sans nécessiter les équipements de lithographie EUV ultra-avancés contrôlés par l'Occident.

Comparatif Calcul Analogique vs Numérique pour l'IA

Critère	Calcul Numérique (GPU)	Calcul Analogique RRAM
Vitesse (matrices IA)	Référence (×1)	×1 000 supérieur
Efficacité énergétique	Référence (×1)	×100 supérieur
Précision	Jusqu'à 32 bits	24 bits (virgule fixe)
Goulot de von Neumann	Oui (CPU ↔ RAM)	Non (compute-in-memory)
Processus de fabrication	TSMC 5nm / EUV requis	Procédés standards (CMOS)
Maturité écosystème logiciel	Très élevée (CUDA)	En développement

Le Contexte Géopolitique — Sanctions et Conséquences Inattendues

La Guerre des Puces : Chronologie d'une Escalade

Depuis 2019, les États-Unis ont progressivement renforcé les restrictions sur les exportations de technologies avancées vers la Chine. Les sanctions d'octobre 2022, sous l'administration Biden, ont marqué un tournant décisif : pour la première fois, Washington cherchait explicitement à freiner le développement technologique d'un rival en interdisant l'exportation des puces IA les plus avancées de Nvidia et AMD, ainsi que des équipements de fabrication de semi-conducteurs de pointe. Les restrictions ont été élargies en 2023 et 2024, ciblant même les versions édulcorées que Nvidia avait développées pour contourner les premières sanctions. La logique était claire : en privant la Chine des outils de calcul les plus performants, on ralentirait sa progression dans les domaines de **l'IA, du supercalcul et des applications militaires**.

L'Effet Boomerang des Sanctions

L'intention des sanctions était limpide, mais leur effet s'est révélé paradoxal. Plutôt que de paralyser l'industrie technologique chinoise, elles ont déclenché une mobilisation sans précédent vers l'autosuffisance en semi-conducteurs. Le gouvernement a injecté des centaines de milliards de dollars dans le secteur. L'histoire de Cambricon Technologies illustre cette dynamique : fondée en 2016, cette startup de puces IA a vu son principal client Huawei l'abandonner en 2019 suite aux sanctions. Ce qui aurait pu être fatal s'est transformé en opportunité — son action a bondi de plus de 765% en 24 mois. Ce phénomène s'est répété avec Huawei (Ascend), Baidu (Kunlun) et de nombreuses startups.

DeepSeek et le Moment Sputnik de l'IA

En janvier 2025, DeepSeek a lancé R1, un modèle open source rivalisant avec les meilleures créations d'OpenAI. La nouvelle a provoqué une onde de choc : Nvidia a perdu 593 milliards de dollars de capitalisation boursière en une seule journée — la plus grande perte de l'histoire boursière américaine. Les médias ont comparé cet événement au *moment Sputnik* de la course spatiale. La percée RRAM s'inscrit dans cette même dynamique : la Chine ne se contente plus de rattraper l'Occident, elle explore des voies alternatives pouvant court-circuiter les avantages technologiques que les sanctions cherchaient à préserver.

Les Applications Potentielles — De la 6G à l'Edge Computing

Communications 6G et Traitement du Signal

Les réseaux 6G utiliseront des techniques avancées comme le **MIMO massif**, où des centaines d'antennes travaillent simultanément. Ces systèmes génèrent des volumes colossaux de données devant être traitées en temps réel — les stations de base actuelles peinent déjà à gérer cette charge avec les processeurs numériques existants. La puce RRAM pourrait transformer cette situation, permettant le traitement de signaux massifs avec une consommation

énergétique minimale. Les tests de l'équipe de Pékin sur la détection de signaux MIMO ont démontré des performances exceptionnelles, surpassant les GPU dans cette tâche spécifique tout en consommant une fraction de leur énergie.

Entraînement et Inférence IA

Former un modèle comme GPT-4 nécessite des milliers de GPU fonctionnant pendant des semaines. La puce RRAM pourrait démocratiser l'accès à l'IA en réduisant drastiquement ces coûts. Les algorithmes d'optimisation de second ordre, particulièrement gourmands en calcul mais plus efficaces pour l'apprentissage, deviendraient viables à grande échelle. Pour l'inférence, une efficacité énergétique centuplée permettrait de déployer des capacités IA dans des environnements auparavant inaccessibles. Ces enjeux sont directement liés aux défis de **l'IA embarquée et de l'inférence locale**.

Edge Computing et Autonomie des Appareils

L'une des implications les plus transformatrices concerne **l'edge computing**. Actuellement, de nombreuses applications IA dépendent de connexions constantes à des serveurs distants. Avec des puces analogiques ultra-efficaces, les smartphones pourraient exécuter localement des modèles de langage sophistiqués, les véhicules autonomes pourraient prendre des décisions critiques sans dépendre du réseau, et les dispositifs médicaux pourraient analyser des données de santé complexes en temps réel. Cette autonomie accrue réduirait également la dépendance aux infrastructures cloud contrôlées par quelques grandes entreprises, renforçant la souveraineté numérique des nations. Ces enjeux rejoignent les thèmes explorés dans notre analyse sur la **réglementation de l'IA et les enjeux éthiques**.

Défis et Perspectives — Le Chemin vers la Commercialisation

Les Obstacles Techniques Restants

Malgré l'enthousiasme, plusieurs défis techniques subsistent. La **scalabilité** reste une question ouverte : les démonstrations concernent des matrices 32×32 à 128×128, pas encore des milliards de paramètres. La **durabilité des cellules RRAM** constitue un autre défi — dégradation après de nombreux cycles d'écriture. La variabilité de fabrication nécessite des techniques de calibration sophistiquées dont l'efficacité à grande échelle reste à démontrer.

L'Écosystème Logiciel : Le Talon d'Achille Potentiel

L'avantage le plus durable de Nvidia réside non dans son matériel, mais dans son écosystème logiciel. **CUDA** bénéficie de décennies de développement et d'une immense communauté : des milliers de bibliothèques optimisées créent un effet de réseau difficile à répliquer. Pour que les puces analogiques atteignent leur plein potentiel, un écosystème comparable devra être développé — compilateurs, outils de débogage, bibliothèques optimisées. Ce défi est considérable mais pas insurmontable : la Chine dispose d'une main-d'œuvre technique massive et a démontré sa capacité à développer des écosystèmes lorsque les circonstances l'exigent.

L'Avantage Stratégique de l'Énergie

Jensen Huang a suggéré que la Chine pourrait gagner la course à l'IA grâce à l'électricité plutôt qu'au silicium. Si les puces chinoises consomment cent fois moins d'énergie pour des performances comparables, l'avantage énergétique devient multiplicatif. La Chine investit massivement dans les énergies renouvelables et dispose de capacités considérables en nucléaire, solaire et éolien. Des gouvernements locaux offrent déjà des subventions pour les centres de données utilisant des puces domestiques.

Implications Géopolitiques et Économiques Mondiales

Vers un Monde Technologique Bipolaire

La percée RRAM accélère la bifurcation de l'écosystème technologique mondial en deux sphères distinctes. D'un côté, les États-Unis et leurs alliés dans des initiatives comme le *Chip 4* (États-Unis, Japon, Taiwan, Corée du Sud). De l'autre, Pékin construit son propre écosystème autonome. Cette fragmentation a des implications profondes pour les entreprises mondiales contraintes de choisir entre deux univers technologiques incompatibles, avec des standards, équipements et écosystèmes logiciels distincts.

Le Réalignement des Chaînes d'Approvisionnement

TSMC produit plus de 90% des puces les plus avancées au monde — une concentration géographique qui représente l'un des risques géopolitiques les plus significatifs. Cette vulnérabilité a catalysé un mouvement de diversification : le CHIPS Act américain avec ses 52,7 milliards de dollars de subventions en est l'illustration. La percée analogique chinoise pourrait modifier ce calcul : si des architectures alternatives contournent les technologies de fabrication dominées par l'Occident, la dépendance aux équipements EUV d'ASML pourrait devenir moins critique pour certaines applications.

L'Enjeu des Matériaux Critiques

La Chine contrôle environ 70% de la production mondiale de terres rares et a déjà utilisé ce levier en restreignant les exportations de gallium et de germanium. Les technologies RRAM pourraient utiliser des matériaux différents des semi-conducteurs traditionnels — modifiant potentiellement l'équation des dépendances. Cette dimension de la rivalité technologique est souvent sous-estimée mais pourrait s'avérer déterminante à long terme.

L'Avenir du Calcul — Hybridation et Convergence

Vers des Architectures Hybrides

L'avenir du calcul sera une hybridation sophistiquée des paradigmes numérique et analogique. Chaque approche excelle dans des domaines spécifiques : le numérique pour la précision absolue, l'analogique pour les calculs parallèles massifs. Les systèmes de demain intégreront des unités analogiques pour les couches d'inférence des réseaux neuronaux, tout en conservant

des processeurs numériques pour la logique de contrôle. Des universités américaines, européennes et asiatiques développent des puces neuromorphiques et des systèmes *in-memory computing* — la percée de Pékin n'est pas un événement isolé mais s'inscrit dans un mouvement plus large.

Le Calcul Quantique : Une Troisième Dimension

Parallèlement au renouveau analogique, le calcul quantique progresse, promettant des capacités révolutionnaires pour certaines classes de problèmes. La Chine investit massivement dans ce domaine avec des réalisations comme le processeur Jiuzhang. L'écosystème de calcul de l'avenir pourrait comprendre **trois piliers complémentaires** : le numérique classique pour la logique générale, l'analogique pour l'IA et le traitement du signal, et le quantique pour les problèmes d'optimisation et de simulation spécialisés.

L'Intelligence Artificielle Comme Enjeu de Civilisation

Au-delà des considérations techniques, la course à l'IA soulève des questions fondamentales sur l'avenir de la civilisation humaine. Les approches chinoise et occidentale de l'IA diffèrent significativement en termes de régulation, de confidentialité et de contrôle gouvernemental. Une bifurcation technologique entraînerait une divergence dans les modèles de société rendus possibles par l'IA — des enjeux qui dépassent largement les préoccupations commerciales ou sécuritaires traditionnelles.

Points clés à retenir

- La puce RRAM de l'Université de Pékin atteint **1 000× les performances GPU H100** avec 100× moins d'énergie grâce au calcul matriciel analogique direct en mémoire.
- Les **sanctions américaines ont paradoxalement catalysé** cette innovation en forçant la Chine à explorer des voies alternatives abandonnées depuis les années 1960.
- La commercialisation est crédible : fabrication sur **procédés standards**, pas de lithographie EUV requise.
- L'**écosystème logiciel** (équivalent CUDA) reste le principal défi ; la Chine a les ressources humaines et le soutien étatique pour le développer.
- L'avenir sera **hybride numérique-analogique-quantique** selon les classes de tâches.
- La bifurcation technologique sino-américaine s'accélère vers deux **écosystèmes incompatibles** avec des implications profondes pour les entreprises mondiales.

Conclusion : Un Nouveau Chapitre de l'Histoire Technologique

La puce analogique de l'Université de Pékin représente bien plus qu'une avancée technique. Elle symbolise un changement profond dans la dynamique de la compétition technologique mondiale. Les sanctions américaines, conçues pour contenir la Chine, ont paradoxalement

catalysé une innovation qui pourrait redéfinir les règles du jeu. En ressuscitant une technologie abandonnée, les chercheurs chinois ont démontré que les plus grandes percées surviennent souvent lorsque des contraintes forcent à explorer des voies alternatives.

Pour les États-Unis et leurs alliés, cette situation appelle à une réévaluation stratégique : investir dans sa propre innovation, y compris dans des paradigmes de calcul alternatifs, pourrait s'avérer plus efficace que de simplement tenter de freiner les concurrents. Pour le monde dans son ensemble, des puces IA plus efficaces pourraient démocratiser l'accès à l'intelligence artificielle, permettant des applications bénéfiques dans la santé, l'éducation et l'environnement. *Dans un éclair de courant électrique traversant des réseaux de mémoire résistive, les mathématiques s'accomplissent à la vitesse de la nature.*

Questions Fréquentes

Qu'est-ce que la technologie RRAM et comment fonctionne-t-elle ?

La RRAM (Resistive Random-Access Memory) est une mémoire non volatile qui encode l'information via des variations de résistance électrique plutôt que des charges. En modifiant la conductivité d'un oxyde métallique entre deux électrodes, elle peut stocker et traiter l'information simultanément, éliminant le goulot d'étranglement de von Neumann. Les cellules RRAM organisées en crossbar arrays effectuent les multiplications matricielles de l'IA directement via les lois d'Ohm et Kirchhoff.

En quoi la puce analogique de Pékin surpasse-t-elle les GPU Nvidia H100 ?

Selon les tests publiés dans *Nature Electronics*, la puce RRAM démontre un débit de calcul plus de 1 000 fois supérieur au GPU H100 pour la résolution d'équations matricielles, avec une efficacité énergétique 100 fois supérieure. L'architecture à double circuit résout le problème centenaire de précision analogique, atteignant 24 bits de précision en virgule fixe.

Comment les sanctions américaines ont-elles contribué à cette innovation ?

Depuis 2022, les États-Unis ont interdit l'exportation des puces IA avancées et des équipements de lithographie EUV. Plutôt que de paralyser l'industrie chinoise, ces restrictions ont déclenché une mobilisation vers des voies alternatives. Privés des GPU numériques, les chercheurs ont revisité le calcul analogique — technologie abandonnée depuis les années 1960 — pour contourner la dépendance aux semi-conducteurs de pointe.

Quelles sont les applications concrètes de cette puce analogique RRAM ?

Les applications prioritaires couvrent : (1) Communications 6G — traitement MIMO massif en temps réel avec une consommation minimale ; (2) Entraînement et inférence IA — réduction des coûts énergétiques des data centers, démocratisant l'accès à la puissance IA ; (3) Edge computing — déploiement de modèles IA sophistiqués directement sur smartphones, véhicules autonomes et dispositifs médicaux, sans connexion cloud.

Quels sont les défis restants avant la commercialisation de masse ?

Trois défis majeurs : la scalabilité (matrices 32×32 à 128×128 démontrées, pas encore des milliards de paramètres), la durabilité des cellules RRAM (dégradation sur cycles répétés), et l'écosystème logiciel. Nvidia a mis des décennies à construire CUDA — développer un équivalent pour l'analogique représente un investissement colossal, que la Chine a les ressources humaines et le soutien étatique pour réaliser.

Sources et références : [ArXiv IA](#) · [Hugging Face Papers](#)

Sources et Références

1. Sun, Z. et al. (2025). *Precise and scalable analogue matrix equation solving using resistive random-access memory chips*. Nature Electronics.
2. South China Morning Post. *China's analogue AI chip could work 1,000 times faster than Nvidia GPU*. Octobre 2025.
3. Bloomberg. *US Sanctions Propel China AI Prodigy to \$23 Billion Fortune*. Novembre 2025.
4. TrendForce. *Chinese Scientists Developed a Novel Chip, Crossing a Century-Old Hurdle*. Octobre 2025.
5. CNBC. *China's strategy in AI race with US — big chip clusters, cheap energy*. Novembre 2025.
6. Nature. *A compute-in-memory chip based on resistive random-access memory*. 2022.

Ayi NEDJIMI Consultants — Expert cybersécurité offensive & intelligence artificielle

ayinedjimi-consultants.fr · ayi@ayinedjimi-consultants.fr

© 2026 — Reproduction interdite sans autorisation.