

أهم ٥ أسباب لتنزيل
قدرات الأحمال التي
يتطلبها عمل الذكاء
الاصطناعي على
IBM Power

أساس موثوق به
لتمكين استراتيجية الذكاء الاصطناعي

جلب الذكاء الاصطناعي Power حيث تتوارد بياناتك

تُعدّ أعباء العمل المدعومة بالذكاء الاصطناعي قيمة، لذا لا تدعها عرضة للتهديدات. توفر خوادم Power حماية لأعباء عمل الذكاء الاصطناعي لديك مع أمان مدمج في كل طبقة من البنية، وتحافظ على الرؤى دون التأثير في الأداء بفضل تشفير الذاكرة الشفاف. عند التعامل مع المهام المعقدة، مثل الذكاء الاصطناعي التوليدى، يمكنك التوسع في الاستدلال بثقة دون القلق بشأن الأداء. مع خوادم Power، ستحصل أيضًا على كشف تهديدات برامج الفدية في أقل من دقيقة،

كلّك بعطل Power Cyber vault™ IBM®، ومدّه التشغيل يصل إلى 99.9999%² لألعاب عمل
كاء الاصطناعي الدائمة والمرنة والمستعدة لمواجهة أي تحدي.

قل من دقيقة %99.9999



توفر خوادم Power ببرمجيات مؤسسيه محسنه بالكامل لتجربه سحابه هجينه سلسه، ما يتيح لك نقل أعباء العمل بين الأنظمة المحلية و IBM Power Virtual Serverg بسهولة. سواء أكنت تعمل على تدريب النماذج في السحابة أم تنفيذ الاستدلال محلياً، تمنحك Power المرونة للقيام بالمزيد. علاوة على ما سبق، تدعم هذه المنصة أكثر من 130 من أدوات وحِزم الذكاء الاصطناعي مفتوحة المصدر، مما يسهل على الفرق لديك بناء النماذج ونشرها وتوسيعها دون أي عوائق.

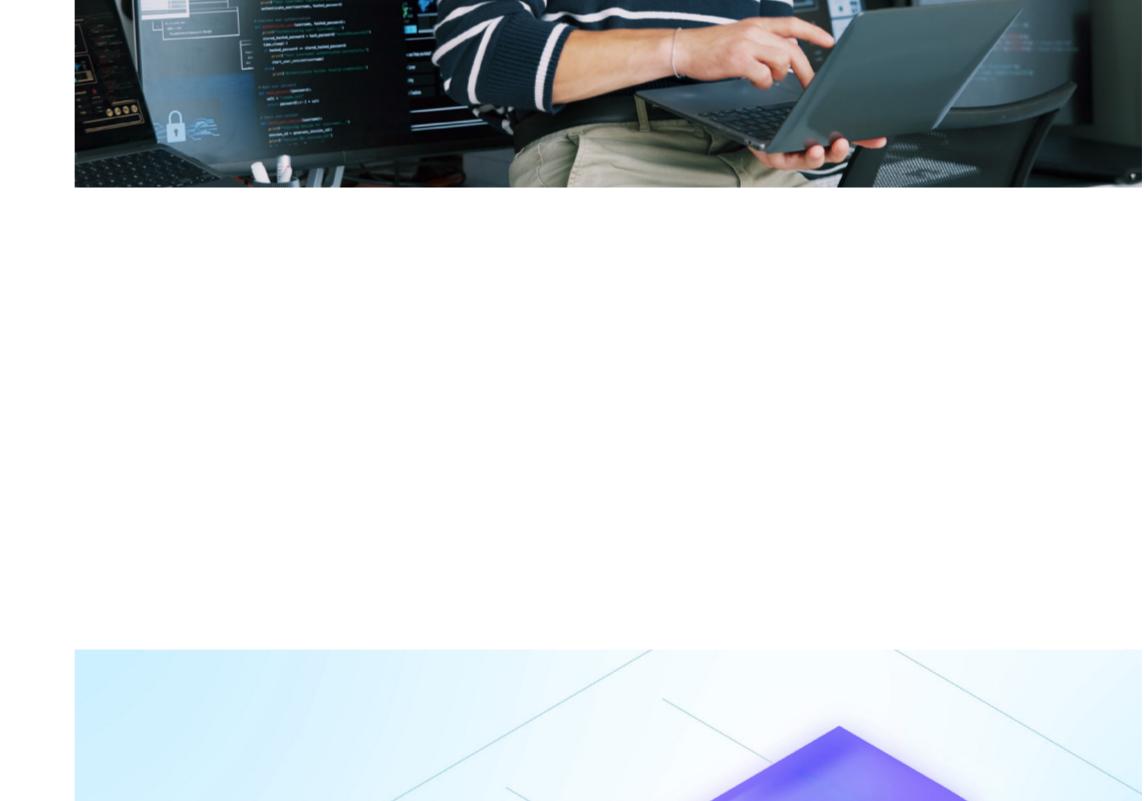
قدّم Power أداءً مستدامًا
عند الطلب.

قال إن متطلبات الاستدامة لا يمكن تحقيقها مع أعباء عمل الذكاء الاصطناعي؟ تقدّم Power11 أداءً أعلى بمرتين لكل واط مقارنةً بالخوادم المعتمدة على معالجات x86، ما يحلك تشغيل أعباء العمل نفسها باستهلاك طاقة أقل.³ وفي ظل الوضع الجديد Energy Efficiency (كفاءة استهلاك الطاقة) في خوادم Power11، يمكنك تحقيق كفاءة طاقة أفضل يصل إلى 28% مقارنةً بالوضع الأقصى.⁴

على نسبة 28%

يسين في كفاءة استهلاك الطاقة مقارنةً بخوادم x86⁴

ووفر Power تسريغاً دون
تقديم تنازلات



الاصطناعي مباشرةً في سير العمل. ولا يحتاج علماء البيانات لديك إلى إعادة البرمجية لاستخدام المنصة؛ إذ يمكنهم تشغيل أعباء عمل الذكاء الاصطناعي إلى ذلك، عند تشغيل نماذج اللغة الكبيرة، يمكنك معالجة عدد أكبر من طلبات في الثانية⁵ على Power S1022 -مقارنةً بالخوادم المعتمدة على تشغيل 40 مستخدماً متزامناً- مع الحفاظ على زمن استدلال يقل عن ثانية رؤى أسرع وعمليات أكثر سلاسة وسرعة استدلال في أقل من الثانية.

%42 ↑

قارنة على إجمالي الإنتاجية بعدد الاستدلالات في الثانية على IBM Power S1022 (20x1 نواة/512 جيجابايت) مع تشغيل 8 SMT مقابل 8 Intel Xeon Platinum 8468V (48x1 نواة/512 جيجابايت). تم إجراء الاختبار باستخدام بيئات Anaconda Python و PyTorch 2.0 و Python 3.10 و Torch 0.9.0. التكوين: حجم الدفع = 60 مع 40 مستخدماً متزامناً. تم تحسين torch.set_num_threads(int) عبر مستويات تحميل Intel Power S1022 لتتمكن من العمل بكفاءة على كلتا المنصتين.

6.26 استدلالات في الثانية مع 40 مستخدماً متزامناً (IBM Power S1022): (<https://www.redbooks.ibm.com/abstracts/redp5675.html>)

النماذج تم ضبطها بدقة بواسطة IBM على مجموعة بيانات داخلية.

كما استندت النتائج إلى اختبار داخلي آخر لأداء الاستدلال على أسئلة وأجوبة باستخدام نماذج PrimeQA (استناداً إلى نماذج Dr. Decr)، صالحة اعتباراً من 31 أغسطس 2023، وتم تنفيذها في مختبر IBM، مع إمكانية اختلاف النتائج حسب حجم أعباء العمل (ColBERT)، باستخدام وحدات التخزين الفرعية وظروف أخرى. النتائج مأخوذة من IBM Power S1022 (20x2 نواة/4-2.9 جيجاهرتز/512 جيجابايت) باستخدام LPAR محاذير NUMA للرقاقة مع 10 نوى. تم إجراء الاختبارات باستخدام بيئات Anaconda Python وبما في ذلك حزم PyTorch 2.0 و Torch 0.9.0، مع استخدام مكتبات Python التي تم تحسينها خصيصاً لمنصة Power. التكوين: torch.set_num_threads(16)؛ حجم الدفع = 1.

(<https://www.redbooks.ibm.com/abstracts/redp5675.html>) IBM Power S1022

(<https://github.com/primeqa>) PrimeQA

