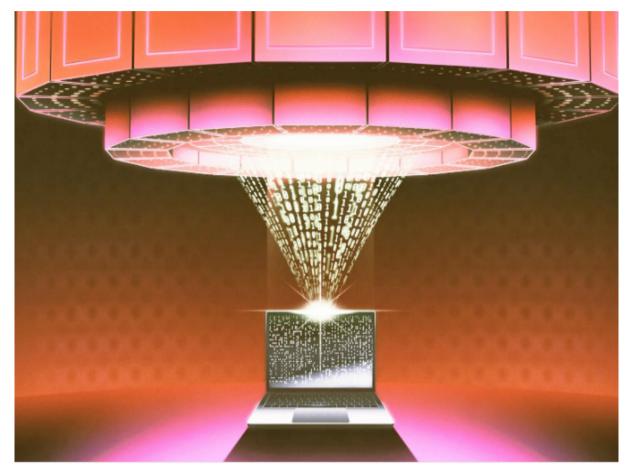
Distillation Can Make Al Models Smaller and Cheaper



THE ORIGINAL VERSION of this story appeared in Quanta Magazine.

The Chinese AI company DeepSeek released a chatbot earlier this year called R1, which drew a huge amount of attention. Most of it focused on the fact that a relatively small and unknown company said it had built a chatbot that rivaled the performance of those from the world's most

famous AI companies, but using a fraction of the computer power and cost. As a result, the stocks of many Western tech companies plummeted; Nvidia, which sells the chips that run leading AI models, lost more stock value in a single day than any company in history.

Some of that attention involved an element of accusation. Sources alleged that DeepSeek had

more efficient way to build Al.

is a widely used tool in Al, a subject of computer a pizza," Vinyals said. The researchers suspected science research going back a decade and a tool that the ensemble models did contain information that big tech companies use on their own models. about which wrong answers were less bad than "Distillation is one of the most important tools others. Perhaps a smaller "student" model could that companies have today to make models more use the information from the large "teacher" efficient," said Enric Boix-Adsera, a researcher model to more quickly grasp the categories it who studies distillation at the University of was supposed to sort pictures into. Hinton called Pennsylvania's Wharton School.

Dark Knowledge

the idea of distilling that onto a single model."

obtained, without permission, knowledge from The researchers thought they might make progress OpenAl's proprietary o1 model by using a technique by addressing a notable weak point in machineknown as distillation. Much of the news coverage learning algorithms: Wrong answers were all framed this possibility as a shock to the Al industry, considered equally bad, regardless of how wrong implying that DeepSeek had discovered a new, they might be. In an image-classification model, for instance, "confusing a dog with a fox was But distillation, also called knowledge distillation, penalized the same way as confusing a dog with this "dark knowledge," invoking an analogy with cosmological dark matter.

After discussing this possibility with Hinton, Vinyals The idea for distillation began with a 2015 paper developed a way to get the large teacher model to by three researchers at Google, including Geoffrey pass more information about the image categories Hinton, the so-called godfather of AI and a 2024 to a smaller student model. The key was homing Nobel laureate. At the time, researchers often in on "soft targets" in the teacher model-where it ran ensembles of models-"many models glued assigns probabilities to each possibility, rather than together," said Oriol Vinyals, a principal scientist at firm this-or-that answers. One model, for example, Google DeepMind and one of the paper's authors- calculated that there was a 30 percent chance that to improve their performance. "But it was incredibly an image showed a dog, 20 percent that it showed a cumbersome and expensive to run all the models cat, 5 percent that it showed a cow, and 0.5 percent in parallel," Vinyals said. "We were intrigued with that it showed a car. By using these probabilities, the teacher model effectively revealed to the

dogs, cats, cows, and cars more efficiently. A big, now been cited more than 25,000 times. complicated model could be reduced to a leaner Considering that the distillation requires access one with barely any loss of accuracy.

Explosive Growth

training data they fed into neural networks, the to distillation. their size.

make smaller models. In 2018, for instance, Google complicated questions. The lab says its fully open researchers unveiled a powerful language model source Sky-T1 model cost less than \$450 to train, called BERT, which the company soon began using and it achieved similar results to a much larger to help parse billions of web searches. But BERT open source model, "We were genuinely surprised was big and costly to run, so the next year, other by how well distillation worked in this setting," said developers distilled a smaller version sensibly Dacheng Li, a Berkeley doctoral student and conamed DistilBERT, which became widely used student lead of the NovaSky team. "Distillation is a in business and research. Distillation gradually fundamental technique in Al."

student that dogs are quite similar to cats, not so became ubiquitous, and it's now offered as a different from cows, and quite distinct from cars. service by companies such as Google, OpenAI, The researchers found that this information would and Amazon. The original distillation paper, still help the student learn how to identify images of published only on the arxiv.org preprint server, has

to the innards of the teacher model, it's not possible for a third party to sneakily distill data from a closed-source model like OpenAl's o1, as The idea was not an immediate hit. The paper DeepSeek was thought to have done. That said, a was rejected from a conference, and Vinyals, student model could still learn quite a bit from a discouraged, turned to other topics. But distillation teacher model just through prompting the teacher arrived at an important moment. Around this with certain questions and using the answers to time, engineers were discovering that the more train its own models-an almost Socratic approach

more effective those networks became. The size Meanwhile, other researchers continue to find of models soon exploded, as did their capabilities, new applications. In January, the NovaSky lab at but the costs of running them climbed in step with UC Berkeley showed that distillation works well for training chain-of-thought reasoning models, Many researchers turned to distillation as a way to which use multistep "thinking" to better answer