

MIT News

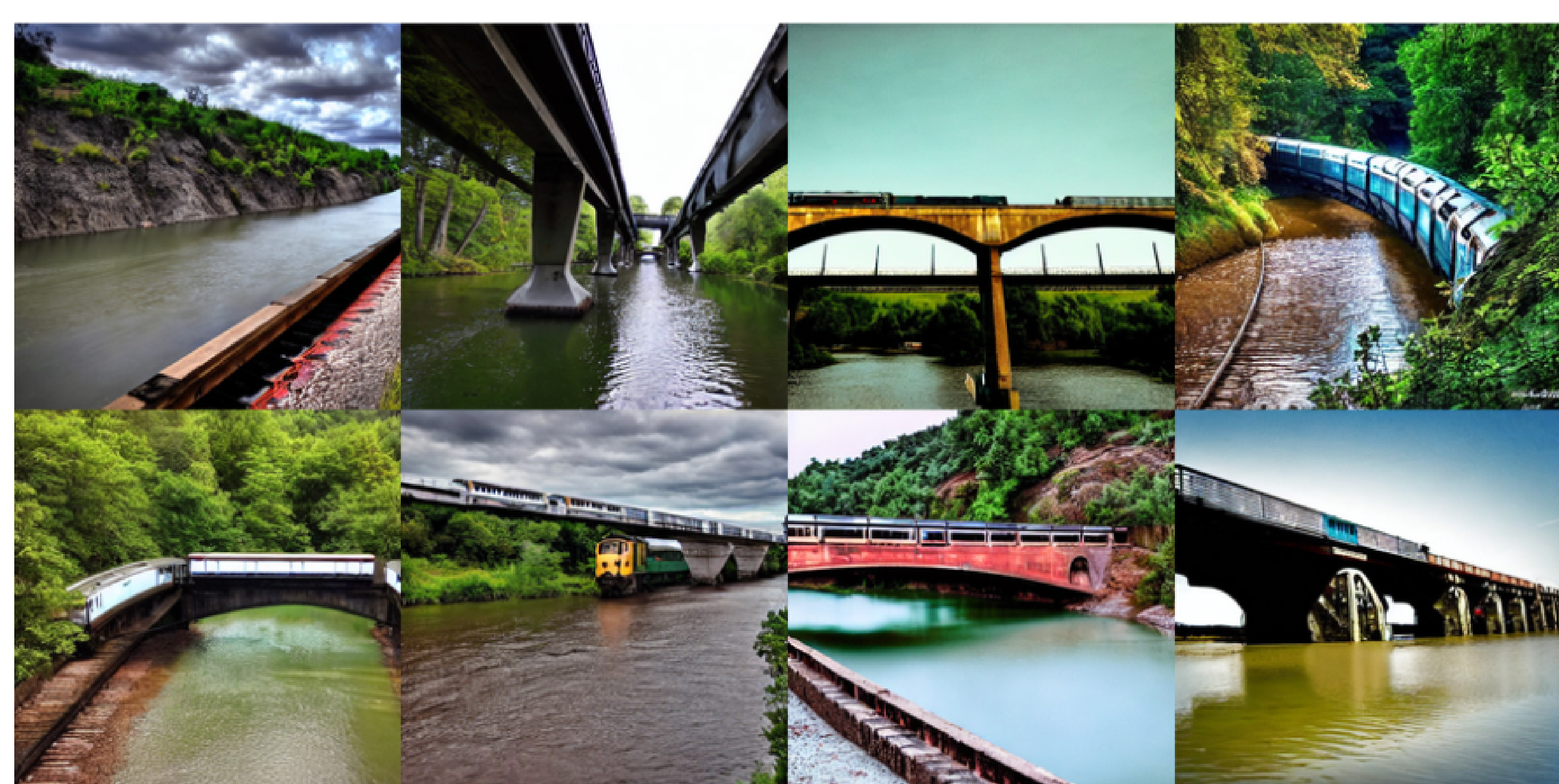
ON CAMPUS AND AROUND THE WORLD

SUBSCRIBE BROWSE SEARCH NEWS

AI system makes models like DALL-E 2 more creative

Researchers develop a new method that uses multiple models to create more complex images with better understanding.

Rachel Gordon | MIT CSAIL
September 8, 2022



This array of generated images on a bridge and "a river under generated using a new method researchers.

Image courtesy of the researchers.

Please answer this nine-question survey to help us make MIT News content as useful and interesting to you as possible. What is your primary reason for visiting MIT News today? Please pick one answer that is the best fit.*

- To read a particular article I saw mentioned somewhere else
- To learn more about MIT
- To find interesting news on science, engineering, or other types of research
- To keep up with news from a particular MIT department, lab, or center
- A different reason (please specify)

NEXT



The internet had a collective feel-good moment with the introduction of DALL-E, an artificial intelligence-based image generator inspired by artist Salvador Dali and the lovable robot WALL-E that uses natural language to produce whatever mysterious and beautiful image your heart desires. Seeing typed-out inputs like "smiling gopher holding an ice cream cone" instantly spring to life clearly resonated with the world.

Getting said smiling gopher and attributes to pop up on your screen is not a small task. DALL-E 2 uses something called a diffusion model, where it tries to encode the entire text into one description to generate an image. But once the text has a lot of more details, it's hard for a single description to capture it all. Moreover, while they're highly flexible, they sometimes struggle to understand the composition of certain concepts, like confusing the attributes or relations between different objects.

To generate more complex images with better understanding, scientists from MIT's Computer Science and Artificial Intelligence Laboratory (CSAIL) structured the typical model from a different angle: they added a series of models together, where they all cooperate to generate desired images capturing multiple different aspects as requested by the input text or labels. To create an image with two components, say, described by two sentences of description, each model would tackle a particular component of the image.

The seemingly magical models behind image generation work by suggesting a series of iterative refinement steps to get to the desired image. It starts with a "bad" picture and then gradually refines it until it becomes the selected image. By composing multiple models together, they jointly refine the appearance at each step, so the result is an image that exhibits all the attributes of each model. By having multiple models cooperate, you can get much more creative combinations in the generated images.

Take, for example, a red truck and a green house. The model will confuse the concepts of red truck and green house when these sentences get very complicated. A typical generator like DALL-E 2 might make a green truck and a red house, so it'll swap these colors around. The team's approach can handle this type of binding of attributes with objects, and especially when there are multiple sets of things, it can handle each object more accurately.

"The model can effectively model object positions and relational descriptions, which is challenging for existing image-generation models. For example, put an object and a cube in a certain position and a sphere in another. DALL-E 2 is good at generating natural images but has difficulty understanding object relations sometimes," says MIT CSAIL PhD student and co-lead author Shuang Li, "Beyond art and creativity, perhaps we could use our model for teaching. If you want to tell a child to put a cube on top of a sphere, and if we say this in language, it might be hard for them to understand. But our model can generate the image and show them."

Making Dali proud

Composable Diffusion — the team's model — uses diffusion models alongside compositional operators to combine text descriptions without further training. The team's approach more accurately captures text details than the original diffusion model, which directly encodes the words as a single long sentence. For example, given "a pink sky" AND "a blue mountain in the horizon" AND "cherry blossoms in front of the mountain," the team's model was able to produce that image exactly, whereas the original diffusion model made the sky blue and everything in front of the mountains pink.

"The fact that our model is composable means that you can learn different portions of the model, one at a time. You can first learn an object on top of another, then learn an object to the right of another, and then learn something left of another," says co-lead author and MIT CSAIL PhD student Yilun Du. "Since we can compose these together, you can imagine that our system enables us to incrementally learn language, relations, or knowledge, which we think is a pretty interesting direction for future work."

While it showed prowess in generating complex, photorealistic images, it still faced challenges since the model was trained on a much smaller dataset than those like DALL-E 2, so there were some objects it simply couldn't capture.

Now that Composable Diffusion can work on top of generative models, such as DALL-E 2, the scientists want to explore continual learning as a potential next step. Given that more is usually added to object relations, they want to see if diffusion models can start to "learn" without forgetting previously learned knowledge — to a place where the model can produce images with both the previous and new knowledge.

"This research proposes a new method for composing concepts in text-to-image generation not by concatenating them to form a prompt, but rather by computing scores with respect to each concept and composing them using conjunction and negation operators," says Mark Chen, co-creator of DALL-E 2 and research scientist at OpenAI. "This is a nice idea that leverages the energy-based interpretation of diffusion models so that old ideas around compositionality using energy-based models can be applied. The approach is also able to make use of classifier-free guidance, and it is surprising to see that it outperforms the GLIDE baseline on various compositional benchmarks and can qualitatively produce very different types of image generations."

"Humans can compose scenes including different elements in a myriad of ways, but this task is challenging for computers," says Bryan Russel, research scientist at Adobe Systems. "This work proposes an elegant formulation that explicitly composes a set of diffusion models to generate an image given a complex natural language prompt."

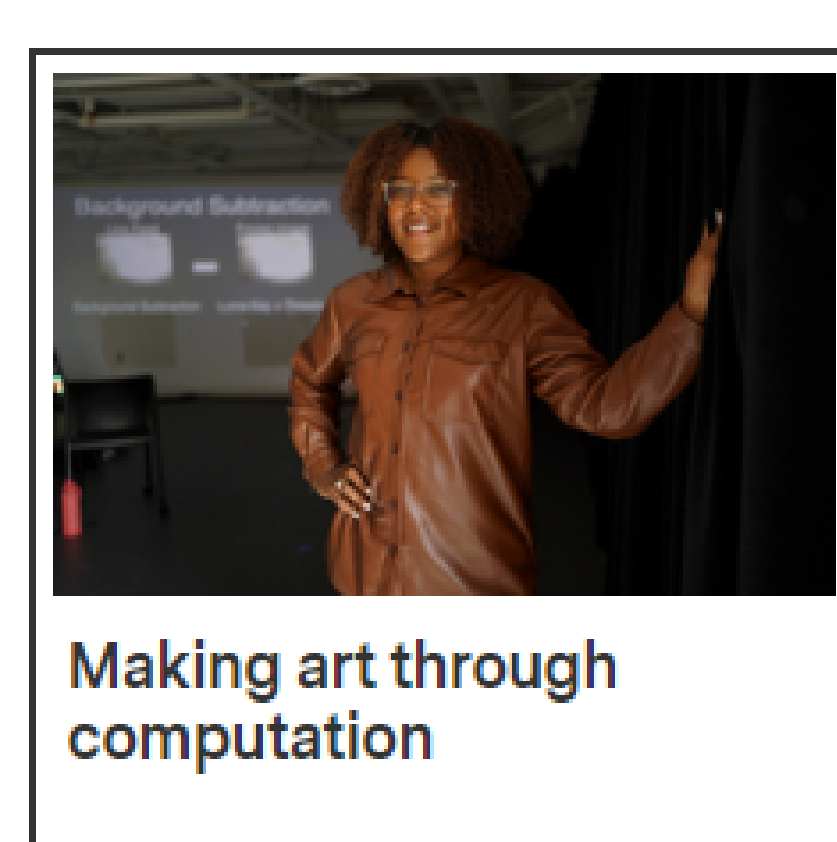
Alongside Li and Du, the paper's co-lead authors are Nan Liu, a master's student in computer science at the University of Illinois at Urbana-Champaign, and MIT professors Antonio Torralba and Joshua B. Tenenbaum. They will present the work at the 2022 European Conference on Computer Vision.

The research was supported by Raytheon BBN Technologies Corp., Mitsubishi Electric Research Laboratory, and DEVCOM Army Research Laboratory.

RELATED TOPICS

- Electrical Engineering & Computer Science (eecs)
- Computer Science and Artificial Intelligence Laboratory (CSAIL)
- Computer science and technology
- Machine learning
- Artificial intelligence
- Algorithms
- Natural language processing
- Computer graphics
- Arts
- School of Engineering
- MIT Schwarzman College of Computing

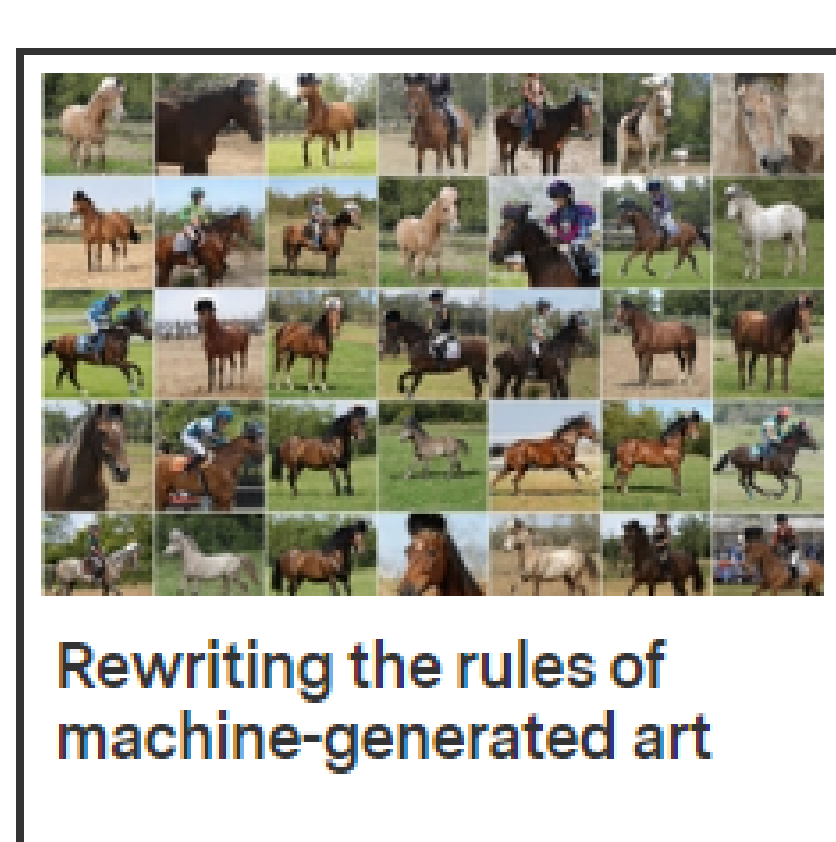
RELATED ARTICLES



Making art through computation



Machines that see the world more like humans do

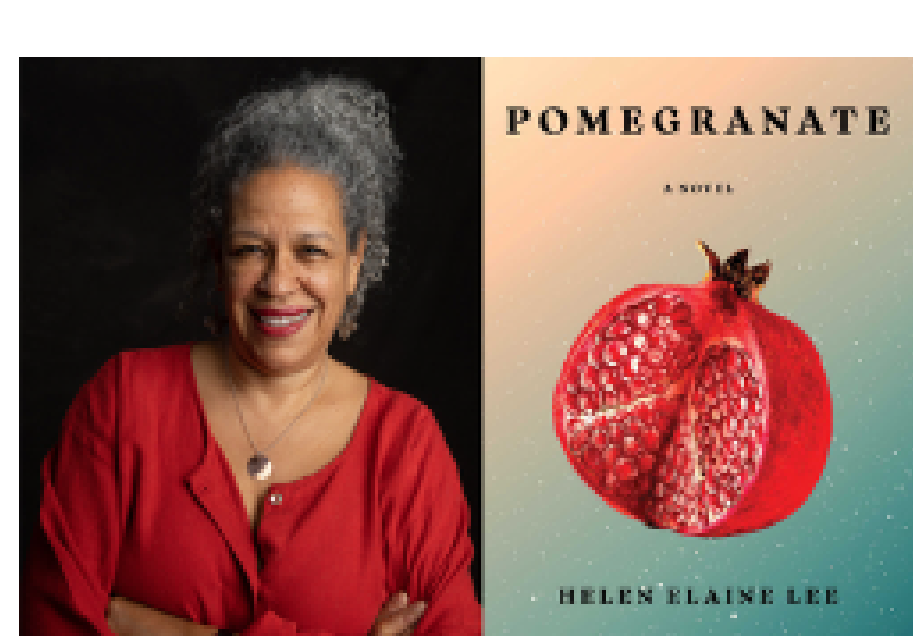


Rewriting the rules of machine-generated art



Educating national security leaders on artificial intelligence

Experts from MIT's School of Engineering, Schwarzman College of Computing, and Sloan Executive Education educate national security leaders in AI fundamentals.



Q&A: A conversation with Helen Elaine Lee about her novel, "Pomegranate"

The MIT professor's new book explores the world of a woman set free from prison and redefining herself in society.



Researchers teach an AI to write better chart captions

A new dataset can help scientists develop automatic systems that generate richer, more descriptive captions for online charts.



Transatlantic connections make the difference for MIT Portugal

The international partnership focuses on climate and sustainability.



Summer 2023 recommended reading from MIT

Enjoy these recent titles from Institute faculty and staff.



Studies at the intersection of equity, computing, and education

"The work I'm doing is deeply rooted in the belief that you can plant seeds in people," says graduate student Cecilié Sadler.

More news on MIT News homepage ->

MIT News

ON CAMPUS AND AROUND THE WORLD

This website is managed by the MIT News Office, part of the Institute Office of Communications.

News by Schools/College:

- School of Architecture and Planning
- School of Engineering
- School of Humanities, Arts, and Social Sciences
- MIT Sloan School of Management
- School of Science
- MIT Schwarzman College of Computing

About the MIT News Office

MIT News Press Center

Terms of Use

Press Inquiries

Filming Guidelines

RSS Feeds

- Subscribe to MIT Daily/Weekly
- Subscribe to press releases
- Submit campus news
- Guidelines for campus news contributors