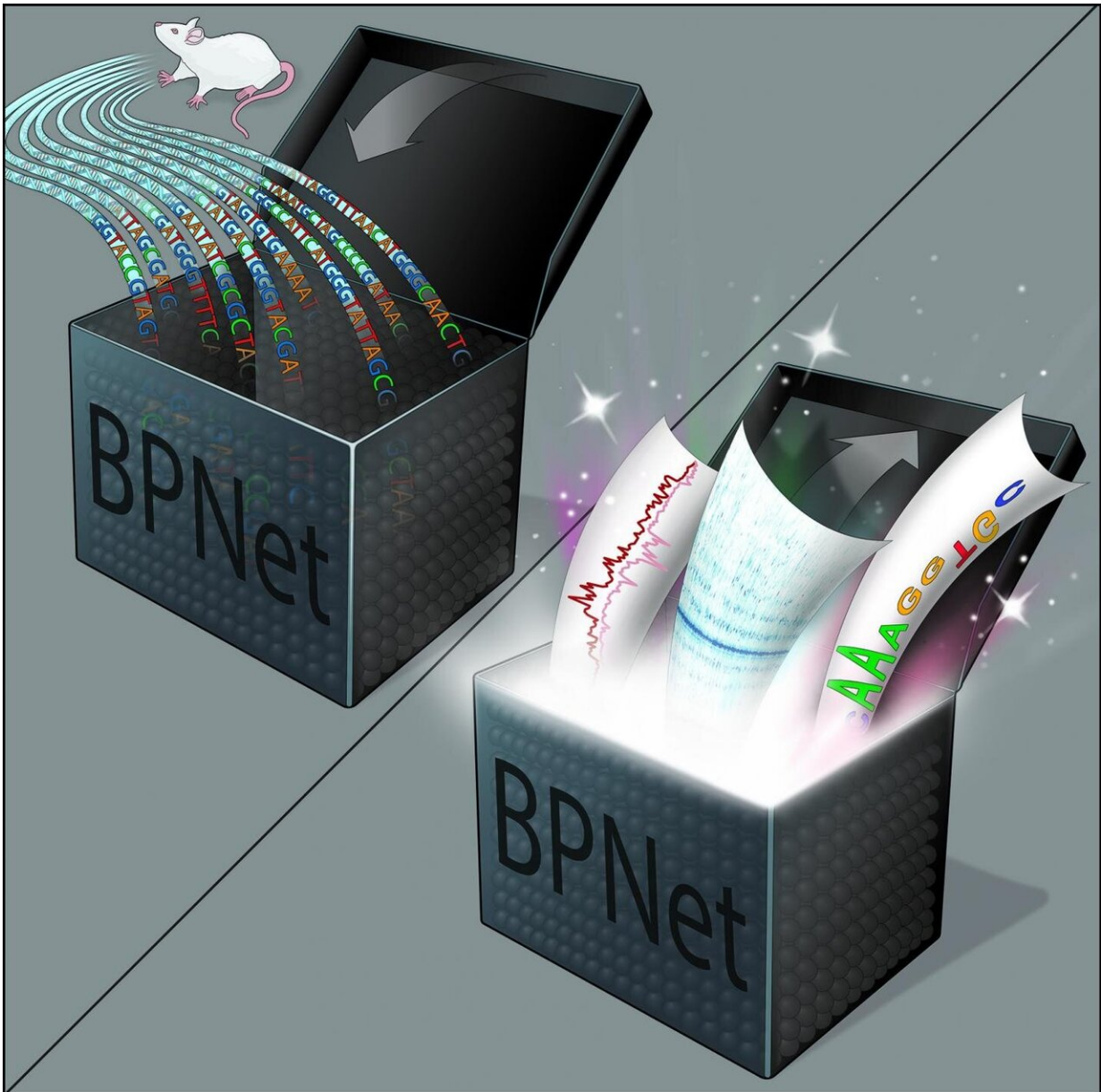


# Explainable AI for decoding genome biology

February 18 2021



Researchers used DNA sequences from high-resolution experiments to train a

neural network called BPNet, whose "black box" innerworkings were then uncovered to reveal sequence patterns and organizing principles of the genome's regulatory code. Credit: Illustration courtesy of Mark Miller, Stowers Institute for Medical Research.

Researchers at the Stowers Institute for Medical Research, in collaboration with colleagues at Stanford University and Technical University of Munich have developed advanced explainable artificial intelligence (AI) in a technical tour de force to decipher regulatory instructions encoded in DNA. In a report published online February 18, 2021, in *Nature Genetics*, the team found that a neural network trained on high-resolution maps of protein-DNA interactions can uncover subtle DNA sequence patterns throughout the genome and provide a deeper understanding of how these sequences are organized to regulate genes.

Neural networks are powerful AI models that can learn complex patterns from diverse types of data such as images, speech signals, or text to predict associated properties with impressive high accuracy. However, many see these models as uninterpretable since the learned predictive patterns are hard to extract from the model. This black-box nature has hindered the wide application of neural networks to biology, where interpretation of predictive patterns is paramount.

One of the big unsolved problems in biology is the genome's second code—its regulatory code. DNA bases (commonly represented by letters A, C, G, and T) encode not only the instructions for how to build proteins, but also when and where to make these proteins in an organism. The regulatory code is read by proteins called [transcription factors](#) that bind to short stretches of DNA called motifs. However, how particular combinations and arrangements of motifs specify regulatory activity is an extremely complex problem that has been hard to pin down.

Now, an interdisciplinary team of biologists and computational researchers led by Stowers Investigator Julia Zeitlinger, Ph.D., and Anshul Kundaje, Ph.D., from Stanford University, have designed a neural network—named BpNet for Base Pair Network—that can be interpreted to reveal regulatory code by predicting transcription factor binding from DNA sequences with unprecedented accuracy. The key was to perform transcription factor-DNA binding experiments and computational modeling at the highest possible resolution, down to the level of individual DNA bases. This increased resolution allowed them to develop new interpretation tools to extract the key elemental sequence patterns such as transcription factor binding motifs and the combinatorial rules by which motifs function together as a regulatory code.

"This was extremely satisfying," says Zeitlinger, "as the results fit beautifully with existing experimental results, and also revealed novel insights that surprised us."

For example, the neural network models enabled the researchers to discover a striking rule that governs binding of the well-studied transcription factor called Nanog. They found that Nanog binds cooperatively to DNA when multiples of its motif are present in a periodic fashion such that they appear on the same side of the spiraling DNA helix.

"There has been a long trail of experimental evidence that such motif periodicity sometimes exists in the regulatory code," Zeitlinger says. "However, the exact circumstances were elusive, and Nanog had not been a suspect. Discovering that Nanog has such a pattern, and seeing additional details of its interactions, was surprising because we did not specifically search for this pattern."

"This is the key advantage of using neural networks for this task," says

Žiga Avsec, Ph.D., first author of the paper. Avsec and Kundaje created the first version of the model when Avsec visited Stanford during his doctoral studies in the lab of Julien Gagneur, Ph.D., at the Technical University in Munich, Germany.

"More traditional bioinformatics approaches model data using pre-defined rigid rules that are based on existing knowledge. However, biology is extremely rich and complicated," says Avsec. "By using neural networks, we can train much more flexible and nuanced models that learn [complex patterns](#) from scratch without previous knowledge, thereby allowing novel discoveries."

BPNet's network architecture is similar to that of [neural networks](#) used for facial recognition in images. For instance, the neural network first detects edges in the pixels, then learns how edges form facial elements like the eye, nose, or mouth, and finally detects how facial elements together form a face. Instead of learning from pixels, BPNet learns from the raw DNA sequence and learns to detect sequence motifs and eventually the higher-order rules by which the elements predict the base-resolution binding data.

Once the model is trained to be highly accurate, the learned patterns are extracted with interpretation tools. The output signal is traced back to the input sequences to reveal sequence motifs. The final step is to use the model as an oracle and systematically query it with specific DNA sequence designs, similar to what one would do to test hypotheses experimentally, to reveal the rules by which sequence motifs function in a combinatorial manner.

"The beauty is that the model can predict way more sequence designs that we could test experimentally," Zeitlinger says. "Furthermore, by predicting the outcome of experimental perturbations, we can identify the experiments that are most informative to validate the model."

Indeed, with the help of CRISPR gene editing techniques, the researchers confirmed experimentally that the [model](#)'s predictions were highly accurate.

Since the approach is flexible and applicable to a variety of different data types and cell types, it promises to lead to a rapidly growing understanding of the regulatory code and how genetic variation impacts gene regulation. Both the Zeitlinger Lab and the Kundaje Lab are already using BpNet to reliably identify binding motifs for other cell types, relate motifs to biophysical parameters, and learn other structural features in the genome such as those associated with DNA packaging. To enable other scientists to use BpNet and adapt it for their own needs, the researchers have made the entire software framework available with documentation and tutorials.

**More information:** Shrikumar, A. et al. Base-resolution models of transcription-factor binding reveal soft motif syntax. *Nat Genet* (2021). [doi.org/10.1038/s41588-021-00782-6](https://doi.org/10.1038/s41588-021-00782-6)

Provided by Stowers Institute for Medical Research

Citation: Explainable AI for decoding genome biology (2021, February 18) retrieved 10 September 2024 from <https://phys.org/news/2021-02-ai-decoding-genome-biology.html>

<p>This document is subject to copyright. Apart from any fair dealing for the purpose of private study or research, no part may be reproduced without the written permission. The content is provided for information purposes only.</p>
--