

# Artificial Intelligence Trends To Watch In 2020

2020

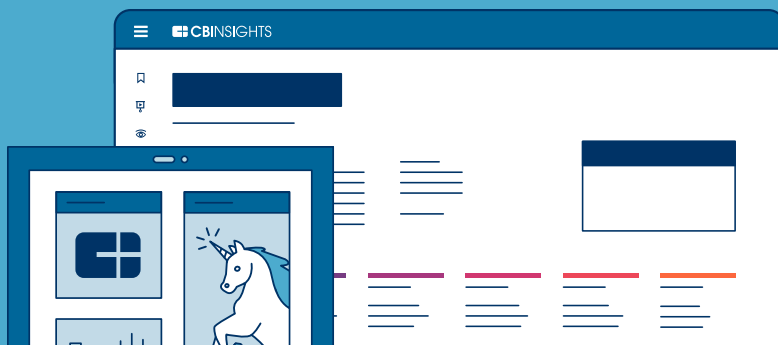


## WHAT IS CB INSIGHTS?

CB Insights is a tech market intelligence platform that analyzes millions of data points on venture capital, startups, patents, partnerships, and news mentions to help you see tomorrow's opportunities, today.

---

[CLICK HERE TO LEARN MORE](#)



## Table of Contents

---

1. Commercial deepfakes will resurrect celebrities, shake up retail, and revolutionize influencer marketing	5
2. Next-gen hacking: AI trojans, voice spoofing, and smart evasion evolve	9
3. AutoML: AI is the future of AI design	13
4. Federated learning will bring in a new data partnership ecosystem	16
5. Alphabet will use AI to dominate smart city contracts	20
6. AI will leave a massive carbon footprint — and we'll need AI to fix it	23
7. Doing more with less: Tackling small data problems in AI will be a major focus	28
8. Quantum machine learning will take baby steps to give traditional AI algorithms a boost	33
9. Natural language processing will help us understand the building blocks of life	37

## AI in 2020

---

Artificial intelligence is fundamentally changing software as we know it.

Corporations are moving past the hype around the technology to discern how it can add practical value. Specifically, data for AI will be a major theme in 2020, as new techniques that train AI with less data or increase data privacy protections gain traction.

Energy efficient AI, quantum neural networks, and the role of natural language processing in understanding amino acids and proteins are some areas where we'll see cutting-edge research this year.

But as the tech matures, it's also bringing along challenges that didn't exist a couple of years ago.

"Deepfakes" are becoming mainstream, making real and fake media indiscernible. AI models can have an inherent bias, and hackers are exploiting the model's bias to fool it. Open-source AI tools — which have been fundamental in democratizing AI — are also easily available for use in malicious applications, like creating the next generation of malware.

In this report, we dive into these and other trends we'll be watching closely in artificial intelligence this year.

# Commercial deepfakes will resurrect celebrities, shake up retail, and revolutionize influencer marketing

---

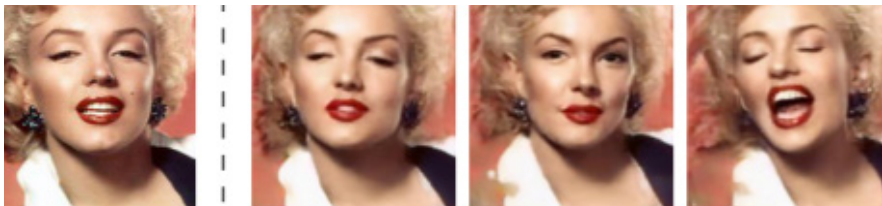
*Deepfakes are hyper-realistic AI-generated images and videos. Media companies are ready to monetize the “benign” side of the controversial tech.*

Deepfakes are controversial, having already made their way into political videos and morphed pornography. But they also provide a great product opportunity for media companies determined to tap into the “fun and goofy” side of the tech.

At the end of December 2019, Snapchat acquired AI Factory, a Ukraine-based startup developing computer vision products, for \$166M.

Snap had previously worked with AI Factory to power Cameos, a feature that enables users to insert selfies into GIFs to create animated deepfakes. Bytedance-owned TikTok is working on a similar feature.

Samsung published a paper on using neural nets to create realistic “talking heads.” Below, the image on the left shows the source, and the ones on the right are AI-generated.



Source: Samsung research

The Financial Times reported on a growing divide between traditional computer generated graphics — which are often expensive and time-consuming — and the recent rise in deepfake tech.

The report says that “deepfakery shaved several years off British actor Bill Nighy in Pokémon Detective Pikachu,” if only for a few frames in the film.

Although a full-fledged, deepfake-driven motion picture may not yet be feasible, Hollywood is heading towards “digitally resurrecting” actors from the '50s and '60s in films — a trend that will benefit from advances in AI.

On the retail front, deepfakes will let brands hyper-personalize visual marketing for consumers. Startup Superpersonal, for instance, swaps out users' faces in short video clips for virtual try-ons.



Source: Superpersonal

**“Consider it you, but as your best fashion model self. The whole process now takes only three minutes (to give an idea of how fast the cloud computing process is moving, back in February it took 20 minutes).”**

**— KATIE BARON, FORBES**

Deepfakes are also making a dent in influencer marketing. Startup [Synthesia](#) used deepfake tech to make David Beckham speak in 9 different languages in a campaign video for the NGO Malaria Must Die. The startup has since raised \$3M from LDV Capital, Mark Cuban, and others.



Source: Malaria Must Die

## IMPLICATIONS

- **Personalization:** Face swapping in retail is all about doubling down on consumer experience. The tech will boost e-commerce experience and virtual online try-ons.
- **Hyper-targeted ads:** As the tech becomes commoditized, hyper-local advertising — such as instant dubbing in different languages — would become low-hanging fruit.
- **Automation in creative fields:** The use of deepfake tech in the TV and film industries could lead to proliferation in sequels, spin-offs, and cultural adaptations of existing content. In other fields, such as casting and modeling, AI-generated faces could stunt demand for human influencers and models.



## Next-gen hacking: AI trojans, voice spoofing, and smart evasion evolve

---

*New-age hacks will evolve on two fronts: fooling AI systems and leveraging AI to launch sophisticated attacks.*

In 2019, researchers from Sydney-based Skylight Cyber found a way to bypass the AI-based antivirus software created by [Cylance](#) — a cyber AI unicorn Blackberry had acquired just a few months prior.

Skylight reported that it found an inherent bias in the AI model and exploited it to create a universal bypass that allowed malware to go undetected.

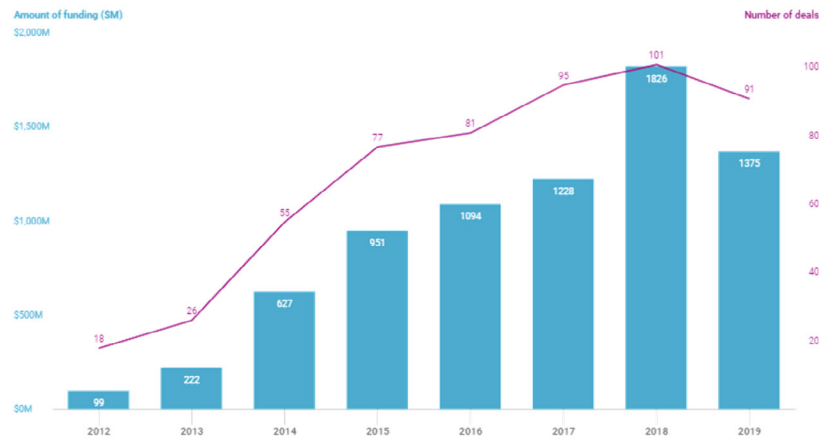
**“...if you could truly understand how a certain model works, and the type of features it uses to reach a decision, you would have the potential to fool it consistently, creating a universal bypass.”**

— SKYLIGHT CYBER

As cyber AI startups raise more funds to protect companies and consumers with the tech, a new crop of hackers and malware that target weaknesses unique to AI will emerge.

## Cyber AI startups will face new attack vectors

Funding to cyber startups that use AI as a core product differentiator



Source: cbinsights.com

 CBINSIGHTS

Hackers can also fool AI through **data poisoning**. Corrupting the training data of AI algorithms impacts how the AI eventually classifies malicious and normal behavior in a network. In other types of adversarial attacks, hackers can introduce small perturbations, which are invisible to the human eye, to an image to trick a neural network into misinterpreting it.

While cyber defenders grapple with these new-age threats that increase with AI adoption, AI itself can be used to create more sophisticated, hyper-targeted cyber attacks.

**Media reports of AI-generated voice spoofing first emerged in Europe in Q3'19.** AI software was used to mimic the voices of CEOs in phone calls to employees with instructions to transfer money, followed by an email with details of the transfer. At least three different companies, including a British energy firm, were targeted.

PRO CYBER NEWS

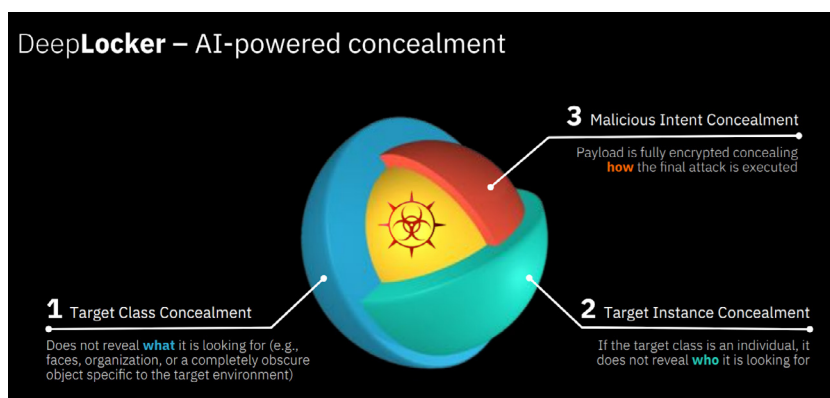
## Fraudsters Used AI to Mimic CEO's Voice in Unusual Cybercrime Case

Scams using artificial intelligence are a new challenge for companies

Source: Wall Street Journal

Deepfakes will likely also be misused in extortion scams.

Although real-world incidents have yet to emerge, IBM developed a proof of concept for a deep learning-powered malware, dubbed DeepLocker, back in 2018. Hackers can make hyper-targeted malware by concealing it within AI code. DeepLocker only unlocks the attack when it encounters a very specific trigger, such as a particular person's face or voice. Otherwise, the code sits undetected in benign everyday applications.



Source: [blackhat USA 2018](#)

With open-source tools available, the barrier for entry into AI is low for hackers.

**“While a class of malware like DeepLocker has not been seen in the wild to date, these AI tools are publicly available... In fact, we would not be surprised if this type of attack were already being deployed.”**

— SECURITY INTELLIGENCE

## IMPLICATIONS

- **Hackers are relentless** and AI tools are more easily available to everyone today than ever before.
- **Cyber AI startups will face new attack vectors.** Hackers have proved it's easy to exploit inherent biases of machine learning models and trick the algorithm.
- **Heavy industries are underprepared.** In the last decade, malware including Stuxnet, BlackEnergy, Havex, Troton, and Industroyer have targeted industrial control systems, from Iranian nuclear plants to Ukrainian power grids. Surveys show asset-heavy industries have not kept up with evolving cyber risks and are not prepared for more advanced threats like AI malware.

## AutoML: AI is the future of AI design

---

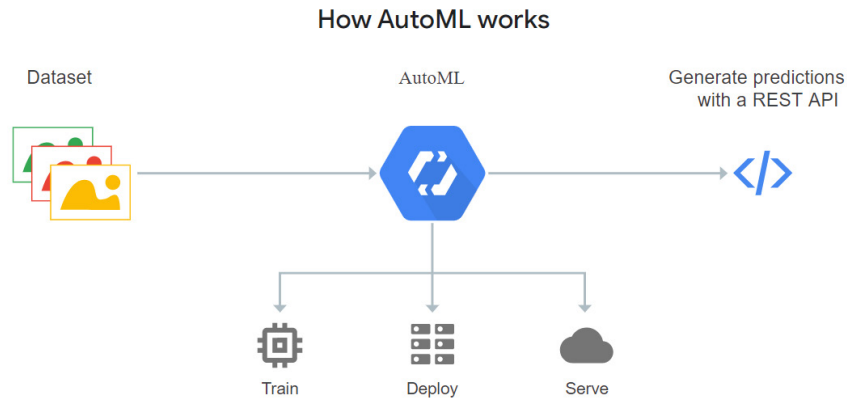
*AutoML, a suite of AI tools to automate design and training of neural networks, will democratize the tech by lowering the barrier to entry for enterprises with minimal AI expertise.*

Designing or searching for the right neural network architecture from thousands of available ones for a specific task is a time-consuming process. It becomes harder when designing AI for more complex situations like autonomous driving, where both speed and accuracy are critical.

Neural architecture search (NAS) is a set of AI approaches to automate the process of finding the best AI design for a given task. Google officially coined the term "AutoML" for this in 2017.

**“Typically, our machine learning models are painstakingly designed by a team of engineers and scientists...if we succeed [with AutoML], we think this can inspire new types of neural nets and make it possible for non-experts to create neural nets tailored to their particular needs...”**

— GOOGLE AI BLOG



Source: Google AI

Since then, adoption of AutoML tools for AI design – including data preparation, training, model search, and feature engineering – has been gradually increasing.

Waymo, for example, recently partnered with Google to automate the process of finding the best neural network architecture to enable autonomous vehicles to identify trees, pedestrians, and vehicles from lidar (light detection and ranging) data.

This adoption is fueled by cloud computing giants like Google making AutoML available to its existing client base. Google Cloud AutoML can be used for computer vision, video processing, translation, and NLP tasks.

Startups are also offering plug-and-play solutions to enterprises. [Databricks](#), a unicorn in the data management and analytics space, introduced AutoML last year. [DataRobot](#), [H2O](#), and [RapidMiner](#) are some other startups that offer AutoML solutions to enterprise clients.

## IMPLICATIONS

AutoML helps to viably scale AI for two reasons:

- **Talent shortage:** There's an acute shortage of AI experts, and AutoML will democratize the tech for enterprises with minimal AI expertise.
- **Cost and complexity:** Designing neural nets is a time-consuming and manual process, even for experts. AutoML can create better solutions and cut down computation costs associated with a trial-and-error approach. We wrote about how AutoML helped Waymo design better AI for perception tasks [here](#).

## Federated learning will bring in a new data partnership ecosystem

---

*Federated learning shows promise in training AI in industries with sensitive and siloed data. In 2020, it will enable a new data partnership model without requiring users to actually share the raw data.*

Federated learning, which allows for increased data privacy while still improving the AI model, is used for applications like Google's text prediction software and URL searches in Firefox.

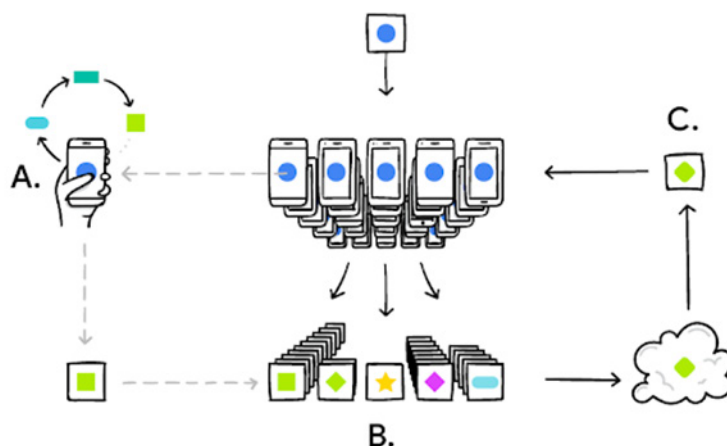
Google initially debuted the tech for its Android keyboard, Gboard, to predict what a user will type next. In the company's Q2'19 earnings call, Sundar Pichai stressed federated learning and other privacy control measures being major focus areas for the tech company.

**“On user privacy and control, it’s always been a big focus for us...Initiative is underway for example like federated learning for almost three years...I think it’s one of the most important areas we are working on.”**

— SUNDAR PICHAI



As depicted below, federated learning allows the Gboard software to improve its AI model without sending raw personal data back to Google.



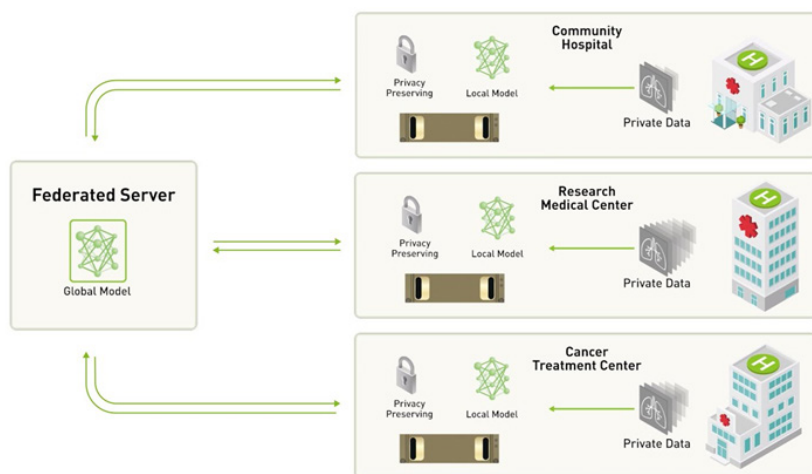
Source: Google

In this case, data stays on your phone instead of being sent to or stored in a central cloud server. A cloud server sends the most updated version of an AI algorithm — called the “global state” of the algorithm — to a random selection of user devices.

Your phone makes improvements and updates to the AI model based on your localized data. Only the update — not the data used to make those updates — is sent back to the cloud to improve the “global state” and the process repeats itself. *(Read about federated learning, and how it differs from other distributed learning approaches [here](#).)*

The ability to protect user data while still improving AI algorithms make federated learning a viable option for industries such as healthcare and banking, where regulatory concerns over data sharing are higher.

For example, Nvidia's AI-powered hardware and software framework for healthcare, called Clara, now supports federated learning. Initial users of the tech include American College of Radiology, MGH and BWH Center for Clinical Data Science, and UCLA Health.



Source: Nvidia

Nvidia also partnered with healthcare startup [Owkin](#), which uses federated learning to predict patients' resistance to cancer drugs. The chipmaker is extending this application to the auto and transportation sectors, enabling cross-country partnerships for accelerating autonomous vehicle research.

In finance, China-based digital bank WeBank is partnering with parent company Tencent's cloud division and Canadian AI research institution Mila for federated learning research.

## IMPLICATIONS

- **Global model, local data:** With federated learning, users can train AI on data stored locally, and only share AI model updates with a cloud server. The “global model” will then benefit all partners in the network to improve their local AI applications.
- **Increase data diversity:** Federated learning will enable more cross-institutional or cross-country partnerships, eventually allowing the global AI model to benefit from diverse local datasets.

# Alphabet will use AI to dominate smart city contracts

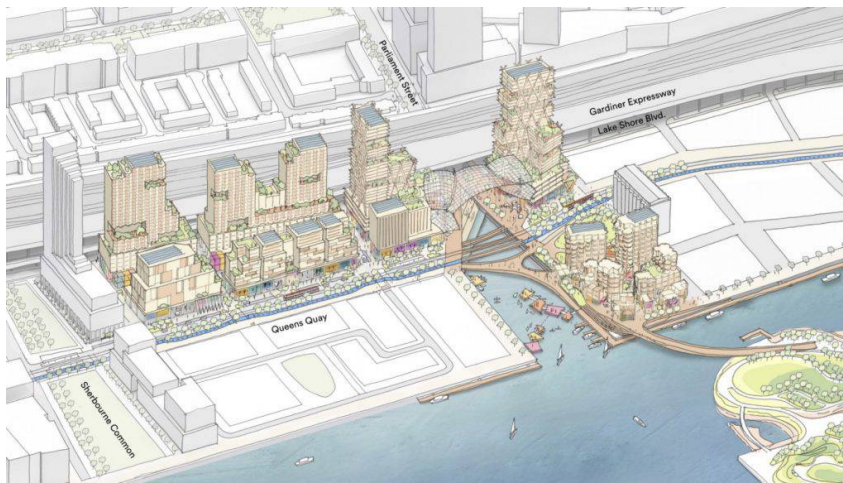
---

*The tech giant is joining forces with local governments to create new city blocks, redrawing competitive lines in sectors ranging from real estate to energy utility to transportation, and more.*

Alphabet, a \$1T AI powerhouse, is making inroads into urban development and smart city planning with IoT and machine learning.

Alphabet subsidiary Sidewalk Labs released a 1,500 page document in Q2'19 highlighting its plans for a \$1.3B smart city development project in Toronto with government and private partners.

This has significant implications for the use of AI in government and urban planning.

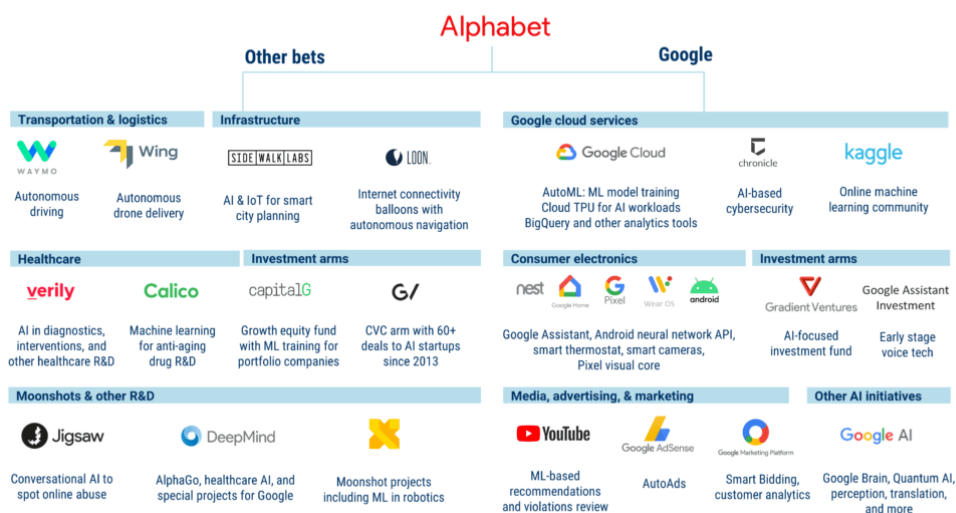


Source: Sidewalk Labs

Smart city planning is a broad concept encompassing smart health, smart mobility, surveillance, and data infrastructure, among other things, with a number of different use cases for AI and machine learning.

Alphabet's reach and AI initiatives in each of those industries make it a formidable competitor for public contractors in sectors ranging from real estate to energy utility to transportation to consulting services.

## Alphabet's org structure: **Key AI initiatives**



For example, Sidewalk Labs has been spinning off smaller companies like Replica and Coord to tackle specific urban development challenges with AI.

[Replica](#) uses machine learning to model commuter behavior and answer questions about what influences a commuter's choices, such as taking public transport. The Portland government will pay Replica over \$450K for a one-year initial service period, and Illinois has signed a \$3.6M deal with the startup for a 3-year period.

Coord uses [machine learning to map curbside assets](#). The company is inviting cities to apply to its 2020 Digital Curb Challenge. Its technology will be offered to winners free of cost, and the pilot will give Coord the opportunity to fine tune its platform and strategies.

Apart from this, Sidewalk Labs' pilot project in Toronto is emphasizing building smart cities with drastically reduced greenhouse emissions and smarter resource management. The company is hiring ML engineers to analyze data from sensors and building management systems to build recommendation engines and predictive models related to energy consumption and sustainability.

## IMPLICATIONS

- **AI expertise makes Alphabet a formidable competitor for government contracts:** Sidewalk Labs has the advantage of pooling in resources and expertise from other Alphabet subsidiaries, including DeepMind, Waymo, and the moonshot lab X.
- **End-to-end solutions:** While a number of vendors offer services like curbside information to cities, Alphabet is well-positioned to offer everything from ML-based urban development tools to autonomous vehicles to building energy management.
- **Financial risk undertaking:** Google is also able to share risk and invest upfront. Sidewalk Labs announced it would bear a "disproportionate share of the cost of upfront innovation" in its public-private partnership model, and receive compensation when it meets performance metrics in later stages. This increases the chances that municipalities and governments will experiment with the tech.

# AI will leave a massive carbon footprint – and we'll need AI to fix it

---

*Given how computationally intensive AI is, the tech will not only necessitate smarter, sustainable solutions – it will also help meet a growing global energy demand.*

One of the reasons many advances in AI have been top-down (i.e., tech giants pioneering AI R&D and open-sourcing tools to others) is how computationally intensive AI research is.

To put this in perspective, Fast Company reported that Google used up “the equivalent of the electricity that the average American household uses in just under six months” for its BigGAN experiment in 2018 to create hyper-realistic images of dogs, butterflies, and burgers.

AI for energy will be a major theme in 2020, from tech giants and automakers to oil & gas companies looking to cut costs, improve efficiencies, and meet global power consumption.

## AI for energy gains media traction

Based on keyword searches “energy” and common AI-related terms



**Two distinct trends are emerging within AI and energy:**

- Energy-efficient AI devices
- AI tools for large-scale energy management

First, energy efficiency will take center stage as AI comes to more edge devices, such as phones and cameras, since edge computing does not have the same power and resources as cloud computing.

For instance, [Kneron](#) is one company that recently announced low-power AI processors for edge devices.

As another example, Apple acquired [Xnor.ai](#), a startup that makes low-power edge AI tools, in Q1'20.

**“...our hardware engineering and machine learning teams asked the audacious question, ‘can we create a hardware, machine learning architecture capable of running deep learning models without a battery? That can be so low-power they can harvest ambient energy from the sun?’”**

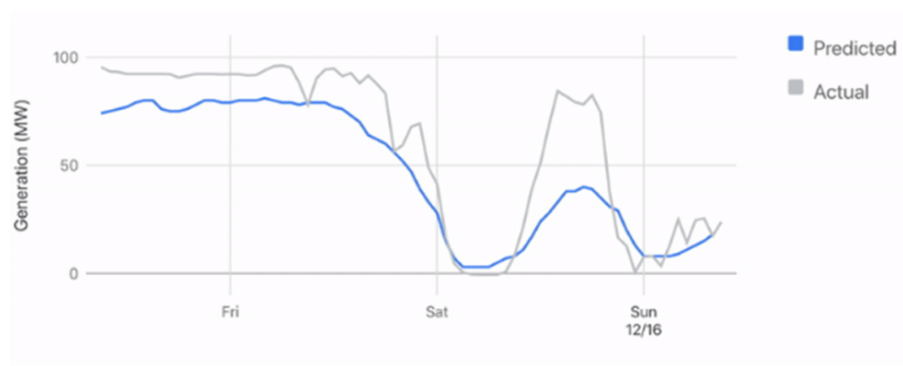
— XNOR.AI



Xnor.ai was working on ultra-low powered cameras that can run AI algorithms. The acquisition was an obvious move for Apple, which is doubling down on AI chips and VR apps for iPhones.

The second trend to watch for in this area is AI-drive energy management and forecasting for large-scale power plants and utilities.

Google has been making a major push towards purchasing 100% renewable energy for its data centers — and it's leveraging AI to help do this. The tech giant partnered with [DeepMind](#) to use neural networks to improve wind energy output.



Source: [Google AI research](#)

DeepMind's neural nets were able to predict wind power output 36 hours in advance based on weather forecasts and wind turbine data.

**“Based on these predictions, our model recommends how to make optimal hourly delivery commitments to the power grid a full day in advance. This is important, because energy sources that can be scheduled (i.e. can deliver a set amount of electricity at a set time) are often more valuable to the grid.”**

— DEEPMIND

Early-stage startups like Bill Gates-backed [Heliogen](#) are already experimenting with niche applications, like using AI to algorithmically control the position of heliostats (mirrors) that capture sunlight.

## IMPLICATIONS

- **Hardware companies will focus on “ultra-low power” devices for ML:** Energy efficiency is a main consideration for AI at the edge (running on devices such as smartphones, smart home cameras, etc).
- **AI for utility-scale energy production:** More cloud giants will transition to using sustainable energy, leveraging AI for increasingly renewable energy output and for streamlining datacenter operations.
- **Streamlining operations for power plants and oil & gas:** Artificial intelligence can predict renewable energy output, automate grid management, aid in precision drilling of oil wells, and power sustainable energy management solutions in smart homes and commercial buildings.

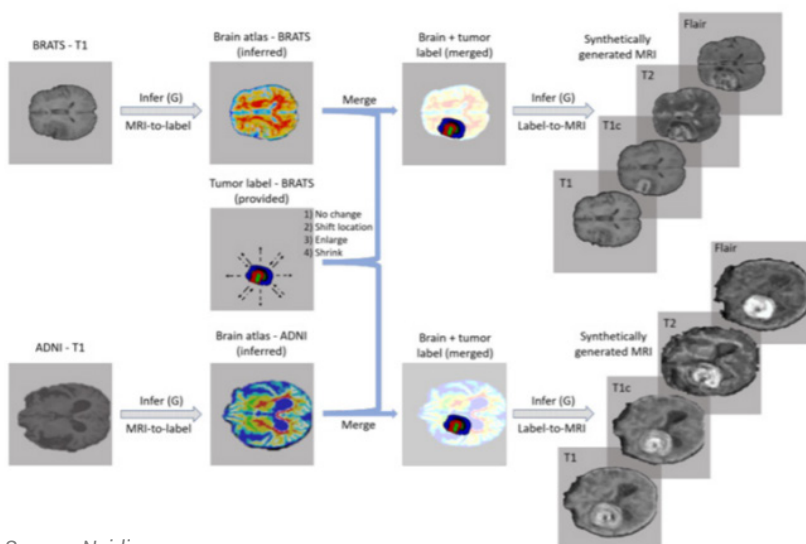
## Doing more with less: Tackling the small data problem in AI will be a major focus

*There are two workarounds if you don't have sufficient data to train data-hungry deep learning algorithms: generate synthetic data or develop AI models that work well with small data.*

Deep learning is extremely data hungry — models are trained on huge sets of labelled data, like millions of images with animals identified and tagged to teach the AI to identify them — and massive amounts labeled data are not available for certain applications. In such cases, training an AI model from scratch is difficult, if not impossible.

One potential solution is to augment real datasets with synthetic data. This has particularly taken off in autonomous driving, where AVs drive millions of miles in photo-realistic simulated environments that can recreate situations like snowstorms and unusual pedestrian behavior, where acquiring real-world data is hard.

Niche synthetic datasets, like the fake MRIs from Nvidia shown below, are also emerging to augment lack sufficient real-world data for rare diseases.



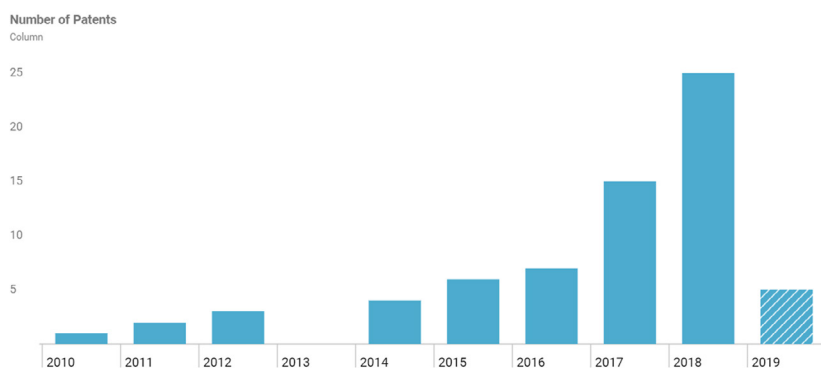
Source: Nvidia

Another way to work around the data issue is to develop AI models that can learn from small data sets.

One approach that has taken off in computer vision tasks is **transfer learning**. This means taking an AI algorithm that's pre-trained for a different task where there is ample labeled data available (for example, identifying cars in images), and transferring that knowledge to a different application for which there is little data (like identifying trucks).

### US patent filings related to 'transfer learning'

One advantage of this AI approach is working around small data problem



Source: cbinsights.com

 CBINSIGHTS

Taking a pre-trained model is like taking a readymade pizza crust and customizing it instead of making the dough from scratch. While it has taken off in computer vision, pre-training was challenging in natural language processing (NLP) given the lack of labeled data in general, until now.

An approach called **self-supervised pre-training** is becoming popular in NLP.

AI is pre-trained on the enormous amount of text readily available on the Internet. For example, [OpenAI](#) pre-trained an AI model with 8M pages of internet text — an extremely compute-intensive task. During training, the model's task was to predict the next word in a sentence based on preceding words.

This is called **self-supervised learning** because there are no “labels” involved here: the AI learns about language by predicting a hidden word based on other words in a sentence.

Researcher Jeremy Howard explains why these self-supervised language models are so important in an excerpt from fast.ai:

**“We are not necessarily interested in the language model itself, but it turns out that the model which can complete this task must learn about the nature of language and even a bit about the world in the process of its training. When we’d then take this pretrained language model, and fine tune it for another task, such as sentiment analysis, it turns out that we can very quickly get state-of-the-art results with very little data.”**

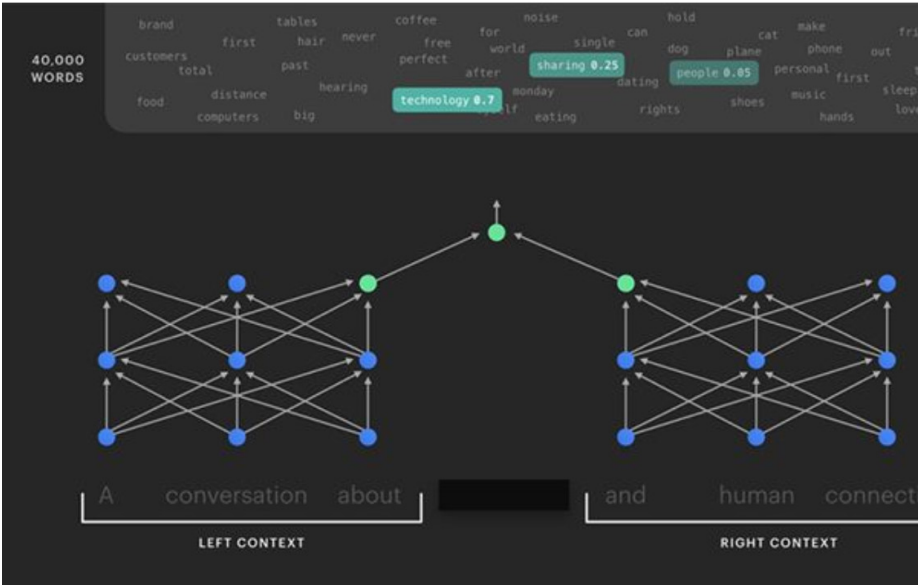
— JEREMY HOWARD, FAST.AI

Another popular example is BERT by Google, where the AI language model not only predicts a word based on the preceding words, but also the succeeding ones (bi-directional understanding of context).

**Input:** The man went to the [MASK]<sub>1</sub> . He bought a [MASK]<sub>2</sub> of milk .  
**Labels:** [MASK]<sub>1</sub> = store; [MASK]<sub>2</sub> = gallon

Source: Google

Facebook's AI division, led by Yann LeCun, has been bullish on self-supervision. One example is pre-training a language model similar to the above examples and fine-tuning it for applications like identifying hate speech.



Source: Facebook

Facebook recently open-sourced its work on self-supervised learning for speech recognition, bypassing the need for manually annotated transcripts by smaller research projects. Facebook open-sourced the code, wav2vec, which is particularly useful for speech recognition in non-English languages, where annotated training data is sparse.

## IMPLICATIONS

- **NLP will hog the limelight in 2020 as a result of self-supervised techniques.** We will finally see better downstream NLP applications like chatbots, advanced machine translation, human-like writing, etc.
- **Big tech players leading the way.** Research on AI models for small data is top-down given how compute intensive it is to develop these pre-trained language models. Tech companies are open-sourcing their research so that other researchers can use it for downstream applications.
- **Synthetic data** and tools that generate realistic fake data level the playing field for smaller companies that don't have access to massive datasets that tech giants do.



# Quantum machine learning will take baby steps to give traditional AI algorithms a boost

---

*Hybrid models, which combine classical machine learning algorithms with quantum AI, will see some practical applications soon.*

Quantum computers require specialized data preparation, quantum algorithms, and quantum programming. In short, the way we interact with classical computers won't work with quantum computers.

Quantum machine learning borrows from the principles of traditional machine learning, but the algorithms are designed to run on quantum processors. This makes them faster than typical neural nets and overcomes hardware constraints that limit current AI research on massive datasets.

**A quick refresher.** Unlike binary computing, where information is stored either as 0 or 1, quantum computers are based on qubits. Qubits can be any value from 0 to 1, or have properties of both of these values simultaneously. Right away, there are [many more possibilities for performing computations](#).

**Despite the hype,** today's most powerful quantum computers, including those being developed by Google, are capable of harnessing the power of 50 to 100 qubits. (To put this in perspective, for quantum computers to have a significant commercial impact, researchers say we require at least a few thousand qubits.)

Research on quantum neural networks (QNN) is very nascent. So given the current hardware limitations of quantum processors, how can QNN algorithms solve any real-world problems?

**“Traditional machine learning took many years from its inception until a general framework for supervised learning was established. We are at the exploratory stage in the design of quantum neural networks.”**

— GOOGLE

Tech giants and quantum startups are looking at a **hybrid approach**, where part of the task is done by traditional neural networks that run on normal computers and another part of it is augmented by QNNs.

Startup [Xanadu](#) applied hybrid classical-quantum AI to transfer learning (see the previous trend on small data for more about transfer learning). The results were promising for image classification tasks.

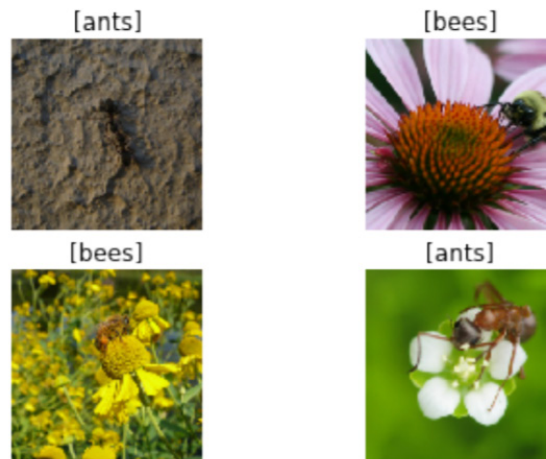


FIG. 3. Random batch of 4 images sampled from the test dataset  $D_B$  and classified by our classical-quantum model (numerically simulated). Predictions are reported in square brackets above each image.

Source: Xanadu research paper, [arxiv.org](https://arxiv.org)

Google's AI team has been focusing on writing algorithms for quantum computers since 2013. The immediate goal, similar to Xanadu's, is to develop "hybrid quantum-classical machine learning techniques on near-term quantum devices."

**"While the current work [on QNN] is primarily theoretical, their structure facilitates implementation and testing on quantum computers in the immediate future."**

— GOOGLE AI BLOG

In two research papers published on the topic, Google explores peculiar ways to train QNNs compared to classical neural net training methods, and tests a QNN's ability to perform simple image classification tasks in simulation.

## IMPLICATIONS

- We will start seeing two of the world's most powerful computing paradigms, quantum computing and AI, solve practical problems initially **in conjunction with classic computers**.
- **Quantum cloud computing** is the latest frontier in cloud wars, with all major providers — AWS, Google, IBM, and Microsoft— doubling down on efforts or venturing into it. This means quantum computers will work in tandem with traditional GPUs and CPUs to add value to cloud clients, as highlighted in a 2020 paper published by Rigetti, Microsoft, and OpenAI. We will see cloud AI algorithms running on such hybrid hardware platforms.

# Natural language processing will help us understand the building blocks of life

*Both natural language processing and genomes are comprised of sequential data. AI algorithms that do well in one field are proving to be useful in the other in surprising ways.*

In the self-supervised learning example discussed earlier in this report, researchers masked specific words in a sentence and had the algorithm guess the missing word to learn about languages more broadly.

Just as sentences are a sequence of words, proteins are a sequence of amino acids in a specific order. Researchers at Facebook AI and NYU used the same concept of self-supervision on massive datasets of protein sequences.

Instead of hidden words, AI has to predict what the hidden amino acid is.

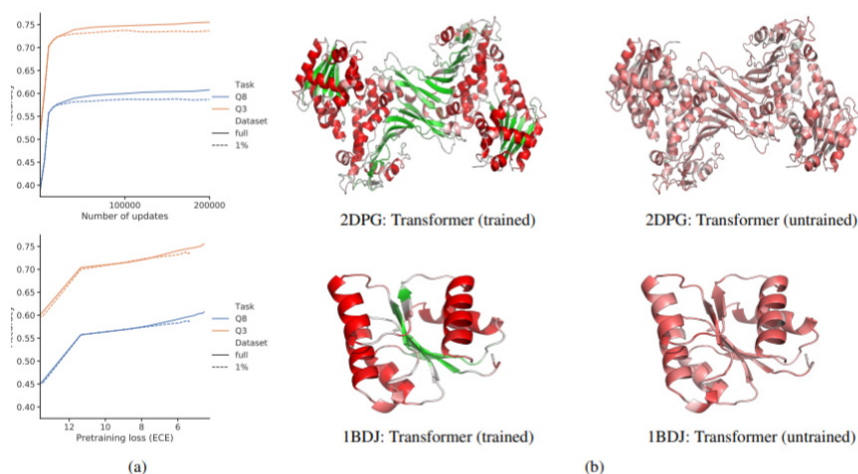


Figure 5: Learned representations encode secondary structure information in a linearly recoverable manner. (a) displays top-1 secondary structure prediction test accuracy of trained Transformer representation projections over the course of pre-training, where top depicts accuracy as a function of model pre-training update steps and bottom depicts accuracy as a function of pre-training loss. (b) illustrates the 3-class representation projections on representative proteins (PDB IDs 2DPG and 1BDJ; Cosgrove et al., 1998; Kato et al., 1999) selected for dissimilarity with the training dataset, having sequence identity no higher than 0.274 and 0.314, respectively, with any training sequence. Red denotes helix, green denotes strand, and white denotes coil. Color intensities denote prediction confidence.

Source: [Biorxiv](#)

Researchers from Germany [tapped into](#) a similar concept of self-supervised language models for classifying proteins.

The most popular recent development was in genomic modeling. **DeepMind** developed an algorithm, **AlphaFold**, to understand protein folding — one of the most complex challenges in genomics — to determine the 3D structure of proteins.

**“...it would take longer than the age of the known universe to randomly enumerate all possible configurations of a typical protein before reaching the true 3D structure — yet proteins themselves fold spontaneously, within milliseconds.”**

#### — DEEPMIND

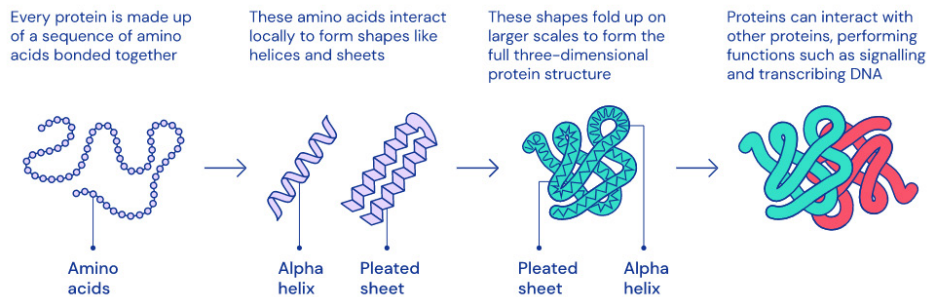


FIGURE 1: COMPLEX 3D SHAPES EMERGE FROM A STRING OF AMINO ACIDS.

Source: [DeepMind](#)

Although AlphaFold uses a hybrid method, it borrows concepts from natural language processing to predict distance and angle between amino acids.

## IMPLICATIONS

- **Better drug design:** Several drug candidates today target proteins, but proteins dynamically change structures based on environmental factors. Understanding their structure and how they fold presents the opportunity to develop drugs for previously unknown targets. Companies like Relay Therapeutics are focused on understanding how proteins move in order to model them, which will aid in new drug discovery.
- AI algorithms can help **model proteins and understand their structure** without requiring in-depth domain knowledge.
- It would be possible to create or optimize **new protein designs** for specific functions in both healthcare and material science.



WHERE IS ALL THIS DATA FROM?

# The CB Insights platform has the underlying data included in this report

---

[CLICK HERE TO SIGN UP FOR FREE](#)

