



Tech Update 2023 – Generative AI Accounting Technology Seminar Series

Randy Johnston

What About Randy?



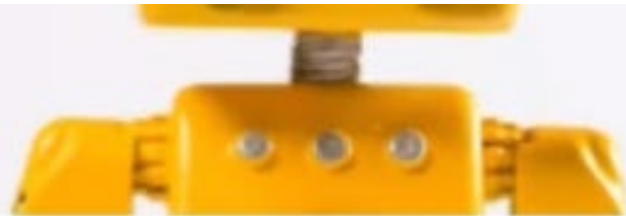
- 40+ years of technology experience, top-rated speaker for almost 40 years
- Top 25 Thought Leaders in Accounting 2011-2023
- 2004-2022 Accounting Today 100 Most Influential in Accounting for nineteen years
- Inducted Accounting Hall of Fame, Feb 2011
- Monthly columns on technology in CPA Practice Advisor
- Published author of six books, From Hutchinson, KS
- randy@k2e.com or randyj@nmgi.com
- 620-664-6000 x 112



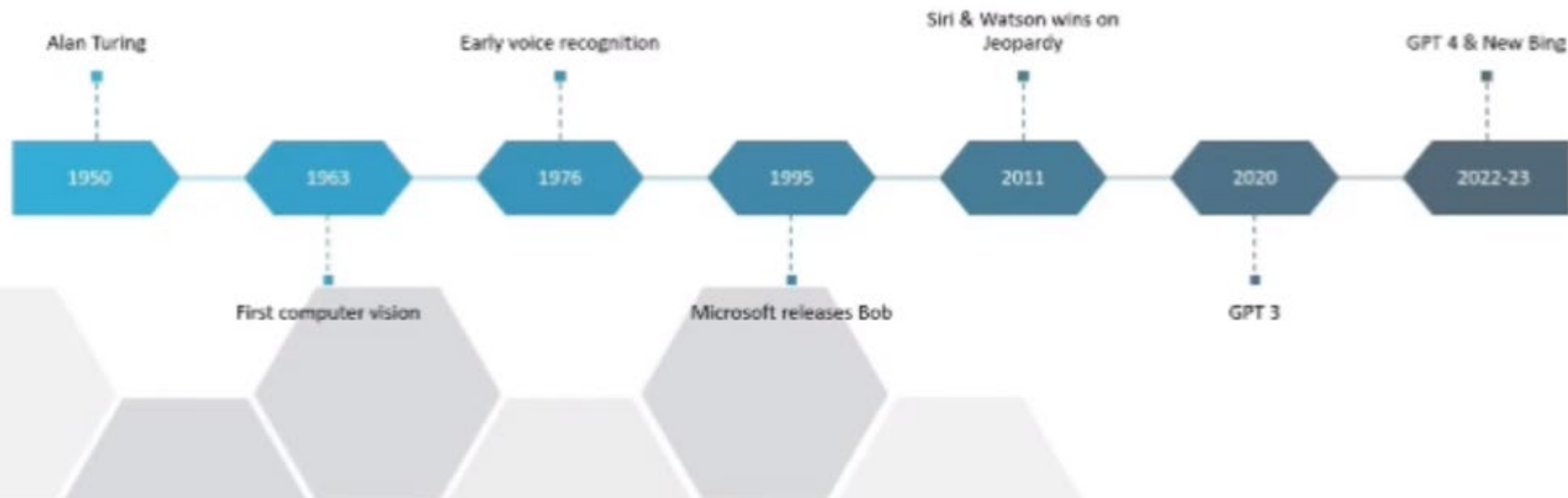


An overview of the technology and how to use it in firms

GENERATIVE AI – CHATGPT IMPACT, LLM, AND REGULATION



EVOLUTION OF AI





Future of Life Institute, An NGO



- In an open letter, “should we automate away all the jobs, including the fulfilling ones? Should we develop non-human minds that might eventually outnumber, outsmart...and replace us? Should we risk loss of control of our civilization?” calling for a six-month “pause” in the creation of the most advanced forms of artificial intelligence (AI) starting with 100’s of signatures and growing to 20,000+
- Italy is the first government to ban ChatGPT as a result of privacy concerns according to the NY Times, 3/31/2023, opened again on May 5

“Large Language Models” (LLMs)



- LLMs—the sort that powers ChatGPT, a chatbot made by OpenAI, a startup—have surprised even their creators with their unexpected talents as they have been scaled up
- Such “emergent” abilities include everything from solving logic puzzles and writing computer code to identifying films from plot summaries written in emoji, passing the bar & USMLE, and has now passed the CPA Exam
- LLMs can, in effect, be trained on the entire internet—which explains their capabilities, good and bad

AI Could Change Computing, Culture, And The Course of History



- As a way of presenting knowledge, LLMs promise to take both the practical and personal side of books further, in some cases abolishing them altogether
- An obvious application of the technology is to turn bodies of knowledge into subject matter for chatbots. Rather than reading a corpus of text, you will question an entity trained on it and get responses based on what the text says. Why turn pages when you can interrogate a work as a whole?
- Bloomberg, a media company, is working on BloombergGPT, a model for financial information. There are early versions of a QuranGPT and a BibleGPT



Uncanny AI



- Such applications and implications call to mind Sigmund Freud's classic essay on the Unheimliche, or uncanny. Freud takes as his starting point the idea that uncanniness stems from "doubts [as to] whether an apparently animate being is really alive; or conversely, whether a lifeless object might not be in fact animate"
- "There's no 'ultimate theoretical reason' why anything like this should work," Stephen Wolfram, a computer scientist and the creator of Wolfram Alpha, a mathematical search engine, recently concluded in a remarkable (and lengthy) blog post trying to explain the models' inner workings



Where Could This Go?



- This raises two linked but mutually exclusive concerns: that AI's have some sort of internal working which scientists cannot yet perceive; or that it is possible to pass as human in the social world without any sort of inner understanding
- “These models are just representations of the distributions of words in texts that can be used to produce more words,” says Emily Bender, a professor at the University of Washington in Seattle
- LLM models are hard to dismiss as “mere babblers”, in the words of Blaise Agüera y Arcas, the leader of a group at Alphabet. Behavior believed to be human is not necessarily so



LLMs Transform Lives and Labor



- ChatGPT embodies more knowledge than any human has ever known
- Contemporary explosion of the capabilities of AI software began in the early 2010s, when a software technique called “deep learning” became popular
- Using the magic mix of vast datasets and powerful computers running neural networks on Graphics Processing Units (GPUs), deep learning dramatically improved computers’ abilities to recognize images, process audio and play games
- But neural networks tended to be embedded in software with broader functionality, like email clients



LLMs-A Giant Exercise in Statistics



- First, the language of the query is converted from words, which neural networks cannot handle, into a representative set of numbers
- GPT-3, which powered an earlier version of ChatGPT, does this by splitting text into chunks of characters, called tokens, which commonly occur together
- These tokens can be words, like “love” or “are”, affixes, like “dis” or “ized”, and punctuation, like “?”
- GPT-3’s dictionary contains details of 50,257 tokens
- GPT-3 vs. GPT-4 has 175 billion parameters vs. 100 trillion



LLMs-Four Key Steps



Tokenization

The promise of large language models is that they _

aptitude talent
potentiality ability
potential capability
promise capacity

Embedding

vocabulary
tongue language
speech

massive
vast huge great
enourmous big
large

facsimile
model replica
imitation duplicate
representation
lookalike

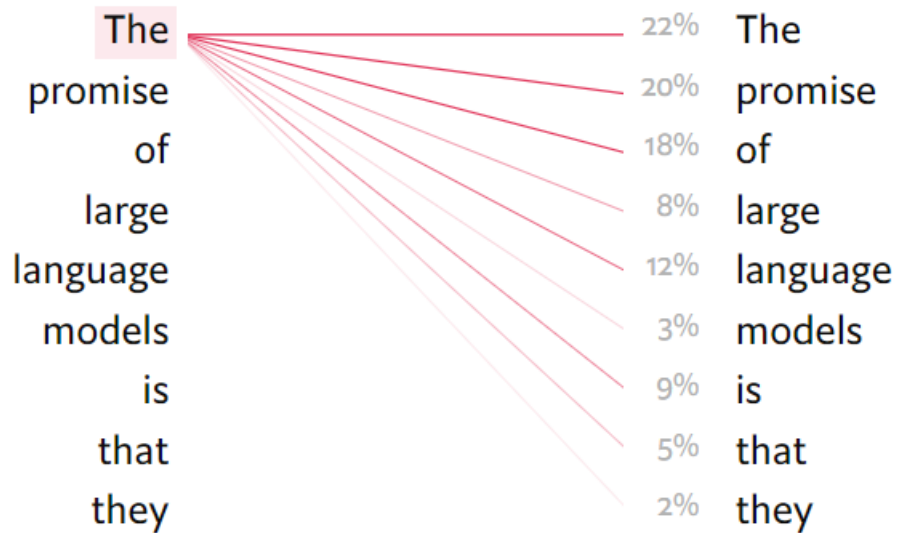




LLMs-Four Key Steps



Attention



Completion

The promise of large language models is that they _

- can 62%
- will 11%
- are 7%
- capture 2%
- could 2%



LLMs-Token Processing



- GPT-3 can process a maximum of 2,048 tokens at a time, which is around the length of a long article
- GPT-4, by contrast, can handle inputs up to 32,000 tokens long—a novella
- The tokens are then **assigned the equivalent of definitions** by placing them into a “meaning space” where words that have similar meanings are located in nearby areas



LLMs-Attention Network



- An LLM then deploys its “attention network” to make connections between different parts of the prompt. Someone reading our prompt, “the promise of large language models is that they...”, would know how English grammar works and understand the concepts behind the words in the sentence. It would be obvious to them which words relate to each other
- An LLM, however, must learn these associations from scratch during its training phase—over billions of training runs, its attention network slowly encodes the structure of the language it sees as numbers (called “weights”) within its neural network

LLMs-Generate Words “Responses”



- Once the prompt has been processed, the LLM initiates a response. At this point, for each of the tokens in the model’s vocabulary, the attention network has produced a probability of that token being the most appropriate one to use next in the sentence it is generating. The token with the highest probability score is not always the one chosen for the response—how the LLM makes this choice depends on how creative the model has been told to be by its operators
- The LLM generates a word and then feeds the result back into itself. The first word is generated based on the prompt alone. The second word is generated by including the first word in the response, then the third word by including the first two generated words, and so on. This process—called autoregression—repeats until the LLM has finished



LLMs-Not Always Predictable



- Jason Wei, a researcher at OpenAI, has counted 137 so-called “emergent” abilities across a variety of different LLMs
- The abilities that emerge are not magic—they are all represented in some form within the LLMs’ training data (or the prompts they are given) but they do not become apparent until the LLMs cross a certain, very large, threshold in their size
- Jonas Degraeve, an engineer at DeepMind, an AI research company owned by Alphabet, has shown that ChatGPT can be convinced to act like the command line terminal of a computer, appearing to compile and run programs accurately. Just a little bigger, goes the thinking, and the models may suddenly be able to do all manner of useful new things

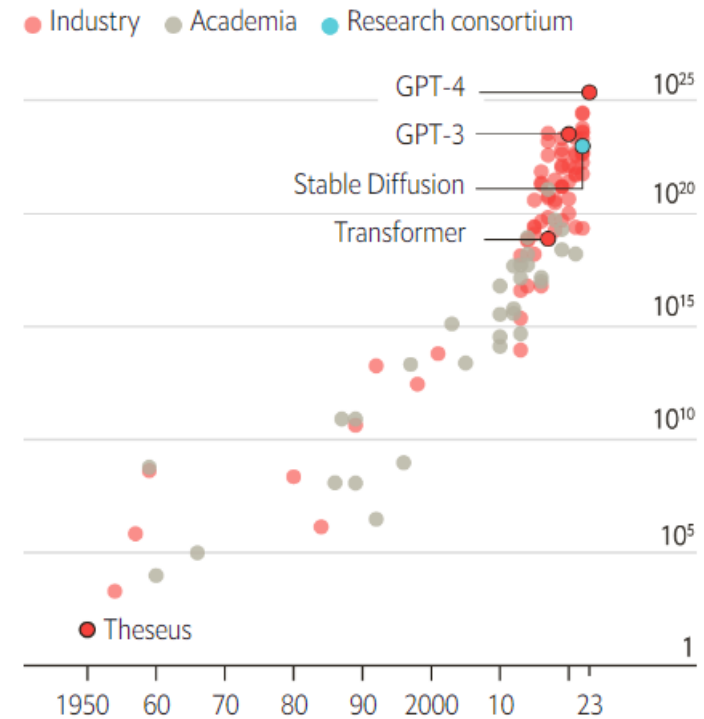
LLMs-Training And Compute



- GPT-3 was trained on several sources of data, but the bulk of it comes from snapshots of the entire internet between 2016 and 2019 taken from a database called Common Crawl of 45TB reduced to 570GB of quality text
- By comparison AlexNet, a neural network that reignited image-processing excitement in the 2010s, was trained on a dataset of 1.2m labelled images, a total of 126 gigabytes—less than a tenth of the size of GPT-4's likely dataset

Faster, higher, more calculations

Computing power used in training AI systems
Selected systems, floating-point operations, log scale



Sources: Sevilla et al., 2023; Our World in Data



LLMs-Attention And Scaling



- “Without attention, the scaling would not be computationally tractable,” says Yoshua Bengio, scientific director of Mila, a prominent AI research institute in Quebec
- GPT-3 has hundreds of layers, billions of weights, and was trained on hundreds of billions of words. By contrast, the first version of GPT, created five years ago, was just one ten-thousandth of the size
- Training GPT-3, for example, used 1.3 gigawatt-hours of electricity (enough to power 121 homes in America for a year), and cost OpenAI an estimated \$4.6m. GPT-4, which is a much larger model, will have cost disproportionately more (in the realm of \$100m) to train

How Generative AI Could Go Wrong

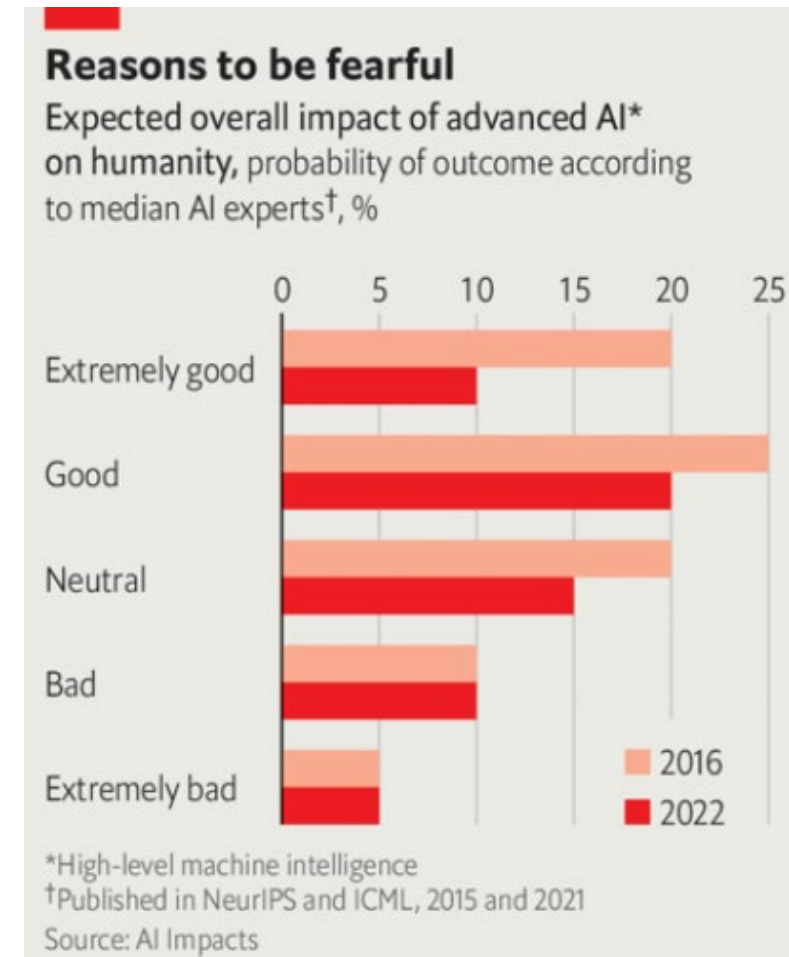


- In 1960 Norbert Wiener published a prescient essay. In it, the father of cybernetics worried about a world in which “machines learn” and “develop unforeseen strategies at rates that baffle their programmers.” Such strategies, he thought, might involve actions that those programmers did not “really desire” and were instead “merely colorful imitation[s] of it.”
- Wiener illustrated his point with the German poet Goethe’s fable, “The Sorcerer’s Apprentice”, in which a trainee magician enchants a broom to fetch water to fill his master’s bath. But the trainee is unable to stop the broom when its task is complete. It eventually brings so much water that it floods the room, having lacked the common sense to know when to stop

How Generative AI Could Go Wrong



- In August 2022, AI Impacts, an American research group, published a survey that asked more than 700 machine-learning researchers about their predictions for both progress in AI and the risks the technology might pose.
- The typical respondent reckoned there was a 5% probability of advanced AI causing an “extremely bad” outcome, such as human extinction (see chart)
- Fei-Fei Li, an AI luminary at Stanford University, talks of a “civilizational moment” for AI
- Asked by an American TV network if AI could wipe out humanity, Geoff Hinton of the University of Toronto, another AI bigwig, replied that it was “not inconceivable”





Human Quality Writing?



- Robert Trager of the Centre for Governance on AI explains, one risk is of such software “making it easier to do lots of things—and thus allowing more people to do them.”
- A text-generation engine that can convincingly imitate a variety of styles is ideal for spreading misinformation, scamming people out of their money or convincing employees to click on dodgy links in emails, infecting their company’s computers with malware
- Chatbots have also been used to cheat at school
- In April, a Pakistani court used GPT-4 to help make a decision on granting bail—it even included a transcript of a conversation with GPT-4 in its judgment



What Are Real Concerns?



- Hallucinations – making things up
- “Alignment problems”, the technical name for the concern raised by Wiener in his essay. An AI might single-mindedly pursue a goal set by a user
- “Reinforcement learning from human feedback” (RLHF). Described in a paper published in 2017, RLHF asks humans to provide feedback on whether a model’s response to a prompt was appropriate
- Another approach, borrowed from war-gaming, is called “red-teaming”. OpenAI worked with the Alignment Research Centre (ARC), a non-profit, to put its model through a battery of tests. The red-teamer’s job was to “attack” the model by getting it to do something it should not, in the hope of anticipating mischief in the real world



Is It Fixable?



- Sam Bowman of New York University and also of Anthropic, an AI firm, thinks that pre-launch screening “is going to get harder as systems get better”
- Another risk is that AI models learn to game the tests, says Holden Karnofsky, an adviser to ARC and former board member of OpenAI
- Another idea is to use AI to police AI. Dr Bowman has written papers on techniques like “Constitutional AI”, in which a secondary AI model is asked to assess whether output from the main model adheres to certain “constitutional principles”



International Agency For AI?



- In the past year alone 37 regulations mentioning AI were passed around the globe
- A poll by the Centre for the Governance of AI found that 91% of a representative sample of 13,000 people across 11 countries agreed that AI needs to be carefully managed
- Call for a global, neutral, non-profit International Agency for AI (IAAI) by Gary Marcus and Anka Reuel
- Interpretability requirements of the AI Bill of Rights proposed by the Biden administration

Review Question



Which of the following is true about artificial intelligence (AI) software?

- a) The products can hallucinate
- b) Client confidential data becomes the property of the publisher
- c) Working with AI is more productive than working without AI
- d) All the above

Review Question



Which of the following is true about artificial intelligence (AI) software?

- a) The products can hallucinate
- b) Client confidential data becomes the property of the publisher
- c) Working with AI is more productive than working without AI
- d) All the above