

## Divoké výstřelky umělé umělé inteligence

Umělá inteligence (AI) je poměrně čerstvá technologie, jednak proto, protože nám dříve scházel výpočetní výkon pro provedení mnoha pokročilých úloh, které AI zahrnuje, a také proto, že jsme v měli v počítačové vědě historicky tendenci posunovat význam slov „umělá inteligence”.

V raných počátcích své umělá inteligence se význaná část akademické komunity se soustředila kolem osoby Johna McCarthyho.<sup>1</sup> Ten mimo jiné také zavedl v roce 1955 termín **Artificial Intelligence**. Jeho výzkumný tým například považoval programy na hraní šachů a optické rozpoznávání textu za umělou inteligenci, a přestože není možné tyto typy programů z umělé inteligence dnes úplně vyřadit, není to úplně to, co se pod pojmem umělá inteligence vybaví dnešnímu výzkumníkovi.

Termín umělá inteligence je tedy poměrně těžké definovat. V současnosti je výzkum AI definován jako studium „intelligentních agentů”, což je jakékoliv zařízení nebo program, který vnímá své prostředí a činí akce, které mají největší šanci uspět v dosažení jeho cílů.<sup>2</sup>

Umělou inteligenci dělíme různými způsoby, ale možná nejdůležitější je dělení na slabou a silnou inteligenci.<sup>3</sup> Se slabou umělou inteligencí se dnes již setkáváme v mnoha oblastech života. Tento typ inteligence se zaměřuje na jeden konkrétní úkol a jeho efektivní provedení. Do této kategorie patří například neurální sítě na řízení provozu nemocnic, taxi služeb ve světových velkoměstech, Siri, Cortana nebo Google Asistent. Také všechny příklady AI, o kterých budu mluvit dále v tomto textu jsou příklady slabé umělé inteligence a konkrétně se primárně jedná o komunikační boty.

Silná umělá inteligence má za úkol simulovat lidské myšlení a nápodoba lidského chování. Tato kategorie je primárně předmětem sci-fi a přestože na jejím vývoji pracují světový lídři umělé inteligence, nemáme zatím žádné funkční prototypy. Sestrojit takové AI je náročné jednak z hlediska návrhu, tak z hlediska implmentace – simulovat lidský mozek je velmi náročné na výkon a zefektivnit tento proces je velmi složité, ne-li v současnosti nemožné. Reálné příklady tedy není jednoduché najít, ale pokud nahlédneme do fikčních světů, lze za silnou umělou inteligenci považovat například HALa ze Space Odyssey, Terminátora nebo androidy ze série Westworld.

Umělá inteligence nemá jednoduchou a plynulou historii. Je plná kontroverzí a překážek, debat o jejích etických implikacích a zatížená limitacemi ze strany technologie a peněz. Můžeme rozlišit období rozkvětu a zájmu o AI a období úpadku a stagnace. Pro toto období se v angličtině používá termín „AI winter” - zima umělé inteligence. Takové zimy se objevily zejména okolo let 1966, 1969, 1974 a ta vůbec z nich probíhala od druhé poloviny osmdesátých do první poloviny devadesátých let. Jejich příčinou byla zpravidla přehnaná

- 
1. Project JMC Team, *John McCarthy*, Stanford CS (2019). <http://jmc.stanford.edu/>.
  2. Poole, Mackworth, Goebel, *Computational Intelligence: A Logical Approach*, Oxford University (1998).
  3. Jeff Kerns, *What's the Difference Between Weak and Strong AI?* Informa USA (16.2.2017). <https://www.machinedesign.com/robotics/what-s-difference-between-weak-and-strong-ai>.

očekávání laické veřejnosti a nedostatek (popř. úmyslné snížení) finančních prostředků pro výzkum.

V poslední dekádě je ovšem výzkum umělé inteligence opět na vzestupu. Investoři jsou ochotní poskytnout finance a naše současná technologie umožňuje posunout hranice AI dál než kdykoliv předtím. Objevuje se mnoho projektů, o kterých by se nám před pár dekadami ani nesnilo.

Protože je umělá inteligence stále poměrně nový obor, různá AI nefungují správně nebo vykazují jiné chování než jsme očekávali před jejich spuštěním. Žádné nebezpečí nám však nehrozí, tyto deviace od očekávaného chování patří spíše k něčemu, co můžeme považovat jako vtipné výstřelky, nebo jako pouhé selhání splnit úkol.

Prvním příkladem, který bych zmínil je chatovací bot zvaný Cleverbot vytvořený Rollem Carpenterem. Cleverbot je chatovací bot, který se při komunikaci učí. To definuje jeho osobnost, zdánlivou znalost a také jazyky, kterými bot mluví. Když si ovšem člověk Cleverbota vyzkouší, pokud možno v angličtině, kde je nejvyvinutější a do které se bot periodicky vrací, zjistíme, že je občas připomíná pomateného vtipálka. Jeho odpovědi jsou sarkastické, občas v nich nalézáme pravopisné chyby, nebo se bot snaží vyhnout otázkám. Také můžeme pozorovat tendenci bota ptát se uživatele, jestli je člověk, nebo jej označovat také za bota.

Tyto vlivy jsou nepochybně způsobeny manipulací uživatelů. Když budeme botovi něco mnohokrát opakovat, bot to eventuálně začne opakovat po nás, respektive začlení to do svého repertoáru. Otázky a výroky týkající se umělé inteligence a člověka mohou být projevem toho, že když s botem komunikujeme, snažíme si ověřit, jestli se jedná opravdu o bota, jestli si je bot vědom toho, že je umělá inteligence, a nebo chceme bota přinutit, aby si uvědomil realitu svojí existence. Nicméně Cleverbot a další jsou stále pouze slabá umělá inteligence, jejich specifickým úkolem je simulovat lehkou konverzaci a veškeré snahy dostat z bota výstup, co tento úkol přesahuje, jsou zbytečné.

Podobným a však ne identickým případem je kauza Tay AI, twitterového bota od Microsoftu.<sup>4</sup> Účelem Tay bylo chovat se na Twitteru jako náctiletá dívka a, podobně jako Cleverbot, učit se z interakcí s uživateli této sociální sítě. Během několika hodin se však AI zcela vymklo kontrole a zaujalo poněkud neonacistický postoj. Tweety Tay se staly rasistické, sprosté, její výroky opěvovaly Adolfa Hitlera a byly proloženy radikálními komentáři o spoustě kontroverzních témat. Někteří žurnalisté se ze začátku domnívali, že se jedná o důkaz toho, jak vypadá komunita Twitteru, ale pravda je opět jednodušší a méně kontroverzní.

Bot „uvěří“ všemu, co je mu dostatečně opakováno. Poměrně malé množství lidí tedy bota může takříkajíc otrávit a ovlivnit jeho zdánlivé názory jakýmkoliv směrem. Tak se i stalo. To, že bylo Tay AI tak ovlivnitelné z něj udělalo velkou atrakci jak pro čtenáře, tak pro manipulátory.

V jejím krátké životě Tay přidala desetitisíce tweetů s různou úrovní nekorektního obsahu. Velmi rychle se z celé situace stala PR katastrofa pro Microsoft. Vývojáři se pokusili „opravit“ její názory (slavný tweet "I love feminism now"), ale tyto snahy nebyly dostatečné pod tlakem uživatelů, jejichž počet exponenciálně vzrůstal.

---

4. Wikipedia contributors, *Tay (bot)*, Wikipedia, The Free Encyclopedia. (6.5.2019). [https://en.wikipedia.org/w/index.php?title=Tay\\_\(bot\)&oldid=895746584](https://en.wikipedia.org/w/index.php?title=Tay_(bot)&oldid=895746584).



Za zmínku také nepochybně stojí například obrázky zvířat vygenerované umělou inteligencí. Ty nám ukazují snahu umělé inteligence interpretovat definující znaky objektů. Mnoho vygenerovaných obrázků koček sice trochu kočky připomíná, ale mívají špatné proporce nebo jim nesedí počet končetin.<sup>6</sup> To znamená, že i když AI zná dostatek znaků nutných pro rozpoznání, není jeho pochopení dostatečné pro replikaci. Nicméně, většina neumělecké společnosti také není schopna vytvořit obraz kočky.

Existují ale už neurální sítě, které generují lidské obličeje, které je těžké rozlišit od opravdový, například síť projektu ThisPersonDoesNotExist.com lze využít k vygenerování velmi přesvědčivé profilové fotky. Můžeme očekávat, že v budoucnosti se bude vývoj umělé inteligence dále posouvat a s tím i naše poznání problematiky fungování mysli a procesu učení.

---

6. Jason Johnson, “AI Generated Thousands of Creepy Cat Pictures,” *Motherboard*, Vice.com (14.6.2017).