

Vyhodnocování inteligence umělých systémů

4IZ431: Přednáška 2

Ondřej Vadinský

KIZI
VŠE Praha

ZS 2021

Obsah

- 1 Úvod
 - Různá pojetí umělé inteligence
- 2 Filosofické předpoklady inteligence
 - Descartovo rozlišení člověka a stroje
 - Turingův test
 - Searlův čínský pokoj
- 3 Kognitivní předpoklady inteligence
 - Úplný Turingův test
- 4 Algoritmická teorie informace
 - Skutečně úplný Turingův test
 - Kognitivní architektury
 - C-Test a Univerzální inteligence
 - (Efektivní) Pragmatická obecná inteligence
 - Kdykoliv přerušitelný test inteligence a Test algoritmického IQ
- 5 Závěr

Úvod

- **Mohou stroje myslet a jak to poznat?**
- *Filosofické a kognitivní předpoklady inteligence.*
- *Přehled přístupů vymezujících, definujících a testujících inteligenci umělých systémů.*
- Definice umělé inteligence: Wang (2019).

Slabá a specifická umělá inteligence

Searle (1980) – slabá UI:

- UI jako nástroj pro poznání lidské mysli,
- program jako přesnější formulace hypotézy o mysli,
- program umožňuje **simulovat mentální schopnosti**.

Goertzel (2014) – specifická UI:

- UI jako tvorba programů pro **řešení specifických úloh**,
- (šachové programy, i sofistikované jako osobní asistentka Siri).

Silná a obecná umělá inteligence

Searle (1980) – silná UI:

- UI (počítač vykonávající správný program) má rozumění a další kognitivní stavy,
- program jako vysvětlení myslí,
- UI je snaha vytvořit **stroj srovnatelný s člověkem**.

Goertzel (2014) – obecná UI:

- UI jako tvorba programů pro **obecné řešení úloh a obecné inteligentní jednání**,
- intelligence je obecná schopnost podkládající jiné schopnosti myslí.

René Descartes

- **Metodologický skepticismus:**
 - Skrze pochybování k pevným principům vědění.
 - „*Myslím, tedy jsem.*“
- **Racionalismus** × empirizmus:
 - *rozum* jako rozhodující zdroj vědění,
 - *smysly* klamou.
- **Dualismus:**
 - nemateriální *mysl*,
 - materiální *tělo a svět*,
 - dobrotivý Bůh zajišťuje shodu při vnímání.

Descartovo rozlišení člověka a stroje I

- Dva předpoklady inteligence:
 - **schopnost rozumné řeči,**
 - univerzálnost myšlení.

Descartes (1637), str. 41:

„...kdyby existovaly stroje, podobající se našim tělům a napodobující naše úkony [...] měli bychom vždy dva velice vážné důvody, abychom poznali, že proto ještě nejsou skutečnými lidmi. První důvod je, že by nikdy nemohly užívat slov ani jiných znaků, skládající je jako činíme my, abychom své myšlenky vyložili jiným. Neboť lze dobře chápat, že stroj může býti udělán tak, aby pronášel slova, ba dokonce aby pronášel některá ve spojení s tělesnými úkony, souvisejícími s nějakými změnami jeho orgánů: jako například když se ho dotkneme na určitém místě, aby se zeptal, co mu chceme říci, když na jiném místě, aby křičel, že ho to bolí, a podobně; nemůže však být udělán tak, aby slova různě sestavoval a takto odpovídal na vše, co se řekne v jeho přítomnosti, jak to i nejtupější lidé mohou činit.“

Descartovo rozlišení člověka a stroje II

- *Dva předpoklady inteligence:*
 - schopnost rozumné řeči,
 - **univerzálnost myšlení.**

Descartes (1637), str. 41:

„...kdyby existovaly stroje, podobající se našim tělům a napodobující naše úkony [...] měli bychom vždy dva velice vážné důvody, abychom poznali, že proto ještě nejsou skutečnými lidmi. [...] A druhý důvod je, že i kdyby vykonávaly určité věci stejně dobře nebo snad i lépe než kdokoli z nás, selhaly by nevyhnutelně v jiných, při nichž by vyšlo najevo, že nejednaly s vědomím, nýbrž toliko podle sestavení svých orgánů; neboť rozum je všestranný nástroj, kterého lze užívat ve všech možných případech, kdežto tyto orgány musí mít nějaké zvláštní uzpůsobení pro každý úkon jednotlivý, a proto je morálně nemožné, aby rozmanitost těchto orgánů v jednom stroji stačila přivést jej k tomu, aby jednal za všech okolností života stejně, jako jednáme my vlivem svého rozumu.“

„Turingův obrat“

Turing (1950):

- **Mohou stroje myslet?**
 - Co je stroj?
 - Co je myslet?
- **Dokáže člověk rozlišit, zda komunikuje s člověkem nebo se strojem?**

Imitační hra I

- *Účastníci:*
 - rozhodčí,
 - muž,
 - žena.
- *Podmínky:*
 - fyzické oddělení účastníků,
 - nepřímá komunikace.
- *Cíl:*
 - **Rozhodnout** pomocí otázek, **kdo je muž** a **kdo žena** (rozhodčí).
 - Zmást rozhodčího (muž).
 - Pomoci rozhodčímu (žena).

Imitační hra II

- *Účastníci:*
 - rozhodčí,
 - člověk,
 - **stroj.**
- *Podmínky:*
 - fyzické oddělení účastníků,
 - **nepřímá komunikace.**
- *Cíl:*
 - **Rozhodnout, kdo je člověk a kdo stroj:**
 - Bude se rozhodčí mýlit stejně často jako v předchozím případě?
 - Bude při 5minutovém rozhovoru úspěšnost běžného rozhodčího nižší než 70 %?

- Digitální počítač (**Turingův stroj**):
 - paměť,
 - výpočetní jednotka,
 - program.
- **Univerzálnost digitálních počítačů**:
 - Turingovy stroje jsou deterministické.
 - Digitální počítač může předpovědět budoucí stavy a tak napodobit libovolný Turingův stroj.
 - Namísto stavby specifického stroje stačí napsat specifický program a ten spustit na libovolném počítači.

Námitky k imitační hře I

Teologická námitka

- *Schopnost myslet vychází z nesmrtelné duše, kterou dal Bůh člověku, ale ne stroji.*
- TURING:
 - Jde proti všemocnosti Boha: Proč by nemohl dát duši stroji?
 - Teologické argumenty bývají nespolehlivé.

„Strkání hlavy do písku“

- *Myslíci stroje? To by bylo strašné! Stroje tedy nemyslí.*
- TURING:
 - Vychází z přesvědčení o nutnosti nadřazenosti člověka.
 - Nejde o solidní argument.

Námitky k imitační hře II

Matematická námitka

- *Existují dokázaná omezení strojů, která neplatí pro lidský rozum:*
 - **Gödelovy věty o neúplnosti** – existují logické formule, které nejsou dokazatelné ani vyvratitelné.
 - **Nerozhodnutelnost problému zastavení** – neexistuje program, který by pro všechny dvojice program – vstup rozhodl, zda se daný program zastaví.
- **TURING:**
 - Nebylo dokázáno, že se omezení netýkají i lidského rozumu.
 - Stroje se tedy mohou v odpovědích mýlit, stejně jako lidé.

Námitky k imitační hře IV

Námitka Lady Lovelace (k Babbageovu analytickému stroji)

- *Analytický stroj není schopen tvořivého myšlení, dělá jen to, k čemu je naprogramován (co programátor ví, jak udělat).*
- *Stroj nás nemůže překvapit.*
- TURING:
 - Nevyplyvá, že nelze naprogramovat myšlení a učení.
 - Stroje nás překvapují, i když je to často chybou v našem očekávání.

Protiargument spojitosti nervového systému

- *Spojitou CNS nelze napodobit diskrétním digitálním počítačem.*
- TURING:
 - Digitální počítač je schopen **imitovat spojitý počítač** dost přesně, aby je rozhodčí nerozeznal.

Námitky k imitační hře V

Neformálnost lidského chování

- *Lidské chování není řízeno sadou formálních pravidel, člověk není stroj.*
- TURING:
 - **Pravidla chování** (rules of conduct) × **zákonitosti chování** (laws of behaviour).
 - Nejsem schopni poznat, zda známe úplnou množinu pravidel.

Námitka extrasenzorického vnímání

- *Telepat by mohl překonat bariéru vnímání, na které je test založen.*
- TURING:
 - Fenomén ač vědecky nepochopený je **statisticky významný**.
 - Zpřísnění podmínek testu („telepathy-proof room“).

- **Stroj je srovnatelný s člověkem:**
 - má mysl,
 - má vědomí.
- *Stroj je lepší než neúspěšný hráč.*
- Odpověď na otázku, zda mohou stroje myslet.

Kritika Turingova testu

- SEARLE:
 - jazyk a rozumění,
 - intencionalita,
- DENNETT:
 - jazyk a svět,
- HARNAD:
 - plný repertoár lidského inteligentního chování,
- SCHWEIZER:
 - externalizmus,
 - inteligentní chování druhu.

Myšlenkový experiment s čínským pokojem II

Searle zamknutý v pokoji:

- **čínské texty:**

- *skript* – znalostní reprezentace,
- *příběh*,
- *otázky*.

- **anglické instrukce:**

- *jak propojit znaky z čínských textů* (na základě jejich tvaru),
- *jak k sadě znaků* (otázka) *sestavit jinou sadu znaků* (odpověď).

- **anglické texty:**

- *příběh*,
- *otázky*.

Myšlenkový experiment s čínským pokojem IV

Důsledky pro silnou umělou inteligenci:

- **Počítače nerozumějí příběhům v čínštině ani jiném jazyce:**
 - protože *Searle jako počítač nerozumí čínsky,*
 - a protože počítače v jiných situacích *nemají nic navíc oproti Searlovi v pokoji.*
- Počítač se správným programem **není postačující ani nutnou podmínkou rozumění:**
 - protože fungující počítač s programem nerozumí,
 - protože *člověk děláající to, co počítač, nezíská rozumění.*
- *Čínský pokoj ukazuje nedostatečnost Turingova testu.*

Úvod	Filosofické předpoklady inteligence	Kognitivní předpoklady inteligence	Algoritmická teorie informace	Závěr	Reference
○ ○	○○○ ○○○○○○○○○○○ ○○○○○●○○○○○	○○○○ ○○○○ ○○○○○○○	○○○○○○○○○○○○○○○ ○○○ ○○○○○○○	○○	

Searlův čínský pokoj

Námitky proti čínskému pokoji I

Námitka systému

- *Searle možná nerozumí čínštině, ale systém jako celek rozumí.*
- SEARLE:
 - Searle je **jediné místo** v systému **schopné rozumění**, a pokud on nerozumí, ani systém nerozumí (*homunkulární strategie*).
 - Searle **se naučí knihy** se znaky a pravidly **zpaměti** a manipulace bude provádět v hlavě, aniž by rozuměl tomu, co dělá (*internalizační strategie*).

Robotická námitka

- *Počítač v robotickém těle bude moci vnímat a jednat a tak rozumět.*
- SEARLE:
 - Nestačí manipulace se symboly, potřebuje i **kauzální vztah ke světu**.
 - Jen **přidá další čínské symboly** na vstupu (vnímání) a výstupu (jednání) a víc instrukcí (program), ale žádné rozumění (*homunkulární strategie*).

Námitky proti čínskému pokoji III

Námitka jiných myslí

- *Rozumění ostatních přisuzujeme jen na základě chování, splní-li počítač behaviorální test, musíme mu také připsat rozumění.*
- SEARLE:
 - Nejde o to, jak vím, ale co přisoudím. To ale musí být víc než jen výpočet.

Námitka mnoha domů

- *Analogové a digitální počítače jsou jen současný stav techniky, v budoucnu půjde sestrojít zařízení s potřebnými kauzálními schopnostmi pro vznik rozumění.*
- SEARLE:
 - Příliš se vzdaluje původnímu vymezení silné UI.
 - Nejde o testovatelnou hypotézu.

Selhání funkcionalizmu

- **Vztah mysl – mozek není jako vztah program – hardware:**
 - Program lze realizovat „*šílenými způsoby*“, kt. zjevně nemají intencionalitu (překládání kamínků).
 - **Program** je čistě formální (**jen syntaxe**), **intencionální stavy** jsou i o obsahu (**sémantika**).
 - Mentální stavy jsou produktem fungování mozku, program není produktem počítače.
- **Simulace není duplikací:**
 - **Počítače nezpracovávají informace, ale manipulují s neinterpretovanými symboly** (interpretace z vnějšku).
 - *Reziduální behaviorismus v UI* (Turingův test).
 - *Reziduální dualismus v UI* (programy jsou zcela nezávislé na hardwaru).

Problémy týkající se mysli a těla

Problém mysli a těla – viz Havel (2001)

- Metafyzický problém: **Zda a jaký je vztah mezi mentálními a tělesnými ději?**
 - *dualismus* – substanční, vlastností, módů bytí,
 - *monizmus* – mentalizmus, idealizmus, materializmus, emergentizmus.

Problém jiných myslí – viz Hyslop (2014)

- Epistemologický a konceptuální problém: **Mají i jiní lidé mysl?**
 - *solipsizmus*,
 - *inference z analogie*,
 - *teorie inferovaná z chování*,
 - *chování jako kritérium přítomnosti mysli*.

Úplný Turingův test I

Harnad (1991):

- *Otázky jako:*
 - „Jsou stroje skutečně inteligentní?“
 - „Rozumějí skutečně stroje tomu, co dělají?“
 - „Mohou stroje skutečně vidět?“
- *jsou vlastně otázkou:*
 - „Mají stroje mysl?“
- **Problém jiných myslí** a jeho řešení je **relevantní pro test UI:**
 - Lidé prisoudí mysl jiným na základě *jejich chování v reálném světě*, které je nerozlišitelné o jejich vlastního chování.
 - *Robot musí v reálném světě* (s lidmi a objekty) **dokázat udělat vše, co běžní lidé**, pro běžného člověka nerozlišitelně od nich.

Úplný Turingův test II

- *Účastníci:*
 - **robot,**
 - **člověk.**
- *Podmínky:*
 - **přímá interakce se světem,**
 - **přímá komunikace.**
- *Cíl:*
 - **Ověřit, zda je robot schopen plného spektra běžného lidského chování.**
 - **Ověřit, zda robot rozumí jazyku:**
 - *jazykový vstup – behaviorální výstup,*
 - *behaviorální vstup – jazykový výstup.*

Externalizmus

Kripke (1972), Putnam (1975) a Burge (1979):

- *Interní reprezentace není dostatečná pro jazykovou referenci:*
 - **přímé kauzální vztahy s prostředím při osvojování jazyka** (intersubjektivně přístupný význam v prostředí),
 - **dělbá lingvistické práce** (experti pro přesné vymezení významu).
- Jazyk je společenský, historicky evolučně vzniklý fenomén.
- **Význam** je veřejná záležitost daná:
 - *pravidelnostmi v mikrostruktuře,*
 - *kauzálními vazbami,*
 - *relevantními experty,*
 - *a zažitou praxí v sociolingvistickém klanu.*

- *Účastníci:*
 - **roboti** (umělý druh),
 - lidé (biologický druh).
- *Podmínky:*
 - přímá interakce se světem,
 - **evoluce komunity.**
- *Cíl:*
 - Ověřit **schopnost robotů** (druhu) **vytvořit si vlastní inteligentní chování,**
 - Ověřit, **zda si roboti vyvinou jazyk.**

Intelligence a kognitivní schopnosti

de Mey (1992):

- **Kognitivní paradigma** – zpracování informací vyžaduje model světa (reprezentaci).
- **Stratifikovaný model vnímání** spojuje:
 - *vjemy objektu,*
 - *očekávání subjektu.*
- **Získávání znalostí** za pomoci vnímání a jednání:
 - jednání nejprve formuje *implicitní znalosti,*
 - opakováním jednání a vnímáním následků dochází k *zexplicitnění implicitního.*

- Různé disciplíny zkoumají kognici na *různých úrovních abstrakce*.
- Vysvětlení kognice pomocí **integrované hierarchie úrovní**:
 - **sociální** – společnost, kultura, organizace, prostředí, ...
 - **psychologické** – individuální chování, zkušenost, dovednosti, ...
 - **komponentové** – kognitivní mechanismy a je realizující komponenty,
 - **fyziologické** – realizace v substrátu.
- Implementace úrovní ovlivněna:
 - *omezeními zdola*,
 - *omezeními shora*.

CLARION

Sun (2007) – kognitivní architektura CLARION:

- *Modulární architektura* (vždy hybridní systém):
 - **akční modul** – řízení vnějších a vnitřních akcí,
 - **neakční modul** – udržování znalostí,
 - **motivační modul** – motivace pro vnímání, jednání a kognici,
 - **metakognitivní modul** – sledování, řízení a úprava ostatních modulů.
- *Požadavky:*
 - učení bez apriorních doménově specifických znalostí,
 - kontinuální učení z interakce se světem,
 - odlišené typy znalostí, včetně specifických metod učení,
 - reaktivní i kognitivní pohled,
 - komplexní situace.

ACT-R

Anderson et al. (2004) – kognitivní architektura ACT-R:

- Základní principy:
 - **racionální analýza** – optimalizace komponent vůči prostředí,
 - nereduktivní rozdělení **znalostí** na **deklarativní** a **procedurální**.
 - **modulární struktura** s komunikací přes zásobníky
- *Skupiny modulů* (různé symbolické systémy):
 - **perceptuálně-motorické** – interakce se světem,
 - **paměťové moduly**:
 - **deklarativní paměť** – fakta jako „chunky“,
 - **procedurální paměť** – „productions“ jak reagovat na situaci.
- *Použití*:
 - paměť, pozornost, přirozený jazyk a komplexní úlohy,
 - neurověda (predikce mozkové aktivity), výuka („cognitive tutors“).

Psychometrická umělá inteligence (PAI)

Bringsjord – Schimanski (2003):

- využívá **psychometrické definice inteligence**.
- Psychometrie – **měření inteligence** a jiných mentálních schopností lidí za použití testů.
- UI by se měla zaměřit na "vytváření takových entit zpracovávajících informace, které budou schopné dosáhnout alespoň **solidních výsledků ve všech zavedených a ověřených testech inteligence a mentálních schopností..**"

Besold et al. (2015):

- **Přímé použití testů lidské inteligence je problematické** (postačující i nutné podmínky).
- Testy potřebují zobecnit a vylepšit.

Univerzální inteligence II

- Abstrahovaná **pracovní definice inteligence**:
„*Intelligence měří schopnost agenta vést si dobře v mnoha různých prostředích.*“
- Formalizace (pokračování):
 - **míra úspěchu v prostředí**:
 - maximalizace očekávaných odměn,
 - suma odměn daných prostředím shora omezená 1,
 - časová preference řešení zabudovaná v prostředí.

$$V_{\mu}^{\pi} := \mathbb{E} \left(\sum_{i=1}^{\infty} r_i \right) \leq 1$$

- **prostor prostředí**:
 - vyčíslitelnost pravděpodobnostní míry prostředí umožňuje testy na počítači.

Univerzální inteligence IV

Legg – Hutter (2007):

- Abstrahovaná **pracovní definice inteligence**:
„*Intelligence měří schopnost agenta vést si dobře v mnoha různých prostředích.*“
- Formální definice **univerzální inteligence**:

$$\Upsilon(\pi) := \sum_{\mu \in E} 2^{-K(\mu)} V_{\mu}^{\pi} \quad \text{kde} \quad V_{\mu}^{\pi} := \mathbb{E} \left(\sum_{i=1}^{\infty} r_i \right) \leq 1$$

Univerzální inteligence různých agentů

- *náhodný* – náhodný výběr akce, neodhalí žádné pravidelnosti v prostředích, typicky nízká V_μ^π , nízká Υ ,
- *specializovaný* – agent velmi dobrý v určité úloze (šachový program), mimo specializaci nízká V_μ^π , nízká Υ ,
- *jednoduchý obecný* – statistika úspěšnosti akcí, náhodný průzkum, odhalí základní pravidelnosti prostředích, vyšší V_μ^π a Υ než předchozí,
- *obecný s historií* – korelace aktuální akce s historií interakcí, odhalí více pravidelností v prostředích, vyšší V_μ^π a Υ než předchozí,
- *obecný s historií a plánováním* – plánování dále než za aktuální odměnu, vyšší V_μ^π a Υ než předchozí,
- *velmi inteligentní agent* – dobré výsledky ve většině jednoduchých prostředí, ale i slušné výsledky v mnoha komplexních prostředích, většinou vysoká V_μ^π ,
- *superinteligence* – vždy vybere akci s nejvyšší budoucí odměnou, dokonalá predikce, maximální V_μ^π ,
- *člověk* – nakolik je omezen evoluční adaptací na určitá prostředí?

Optimální agent AIXI

Legg – Hutter (2007) a Hutter (2012) – rovnice AIXI:

$$a_k := \arg \max_{a_k} \sum_{o_k r_k} \dots \max_{a_m} \sum_{o_m r_m} [r_k + \dots + r_m] \sum_{q: U(q, a_1 \dots a_m) = o_1 r_1 \dots o_m r_m} 2^{-l(q)}$$

„Sequential decision theory (Bellman’s equation) formally solves the problem of rational agents in uncertain worlds if the true environmental probability distribution is known. If the environment is unknown, Bayesians replace the true distribution by a weighted mixture of distributions from some (hypothesis) class. Using the large class of all (semi)measures that are (semi)computable on a Turing machine bears in mind Epicurus, who teaches not to discard any (consistent) hypothesis. In order not to ignore Ockham, who would select the simplest hypothesis, Solomonoff defined a universal prior that assigns high/low prior weight to simple/complex environments, where Kolmogorov quantifies complexity.“

Vlastnosti Univerzální inteligence

- **validní** – definuje inteligenci, ne nějakou související kvalitu, nebo jen aspekt inteligence,
- **smysluplná** – absolutní porovnatelné měřítko (skalární hodnota),
- **široký rozsah** – od velmi nízké po velmi vysokou inteligenci,
- **obecná** – srovnání značně odlišných (nejlépe všech) agentů,
- **nepředpojatá** – testové úlohy nejsou výsledkem kulturní nebo jiné předpojatosti, ale závisí na referenčním Turingově stroji \mathcal{U} ,
- **fundamentální** – založena na „pevných základech“ (koncepty výpočtu, informace a složitosti),
- **formální** – jasně definovaná matematická rovnice,
- **objektivní** – nezávisí na subjektivních kritériích,
- **univerzální** – není antropocentrická,
- **není praktickým testem** – Kolmogorovská složitost je nevyčíslitelná, převod na test bude vyžadovat aproximaci.

Námitky k Univerzální inteligenci II

Předpoklad vyčíslitelnosti prostředí příliš silný

- *Vesmír není vyčíslitelný.*
- LEGG A HUTTER:
 - Vyčíslitelnost ve smyslu vyčíslitelné pravděpodobnostní míry nad budoucími událostmi.
 - Nejsou důkazy pro nevyčíslitelnost ve fyzice v tomto smyslu.
 - Vše dosud známé je v tomto smyslu vyčíslitelné.

Prostředí vracející odměnu s omezeným součtem je nerealistický

- *Vesmír nevrací odměny.*
- LEGG A HUTTER:
 - Jde o testovací framework, subjektivní odměny by neumožnily testovat.

Pragmatická obecná inteligence

Goertzel (2010):

- *Kritika univerzální inteligence:*
 - reální **agenti** jsou **přízpůsobení některým prostředím** – „biased generality“ namísto univerzálnosti,
 - ne všechno inteligentní chování je o vyhledávání odměny,
 - **odměny a cíle** jsou často **dané samotným agentem**,
 - externí odměny jako testovací framework pro dedukci inteligence.
- **Pragmatická obecná inteligence:**
 - lze zvolit pravděpodobnostní distribuci prostředí,
 - podmíněná distribuce cílů v prostředích,
 - vyhodnocování po časový interval od volby cíle do jeho dosažení.

„Intelligence měří schopnost agenta dosahovat komplexních cílů v komplexních prostředích.“

$$\Pi(\pi) \equiv \sum_{\mu \in E, g \in G, T} \nu(\mu) \gamma(g, \mu) V_{\mu, g, T}^{\pi}, \text{ kde } V_{\mu, g, T}^{\pi} \equiv \mathbb{E} \left(\sum_{i=s}^t r_g(I_{g,s,i}) \right)$$

Efektivní pragmatická obecná inteligence

Goertzel (2010):

- *Kritika univerzální inteligence:*
 - reální **agenti** operují s **omezenými zdroji**.
- **Efektivní pragmatická obecná inteligence:**
 - normalizace pragmatické obecné inteligence výpočetními nároky agenta,
 - podmíněná distribuce spotřeby výpočetních zdrojů.

$$\Pi_{\text{Eff}}(\pi) \equiv \sum_{\mu \in E, g \in G, Q, T} \frac{\nu(\mu) \gamma(g, \mu) \eta_{\pi, \mu, g, T}(Q)}{Q} V_{\mu, g, T}^{\pi}$$

Kdykoliv přerušitelný test inteligence

Hernández-Orallo – Dowe (2010):

- *Návrh* testu inteligence pro:
 - současné i budoucí umělé i biologické agenty,
 - různé úrovně inteligence i různé časové škály.
- *Dosažení vyčíslitelnosti Univerzální inteligence:*
 - konečný vzorek prostředí – problém **diskriminační síly prostředí** vyřešen požadavkem **citlivosti vůči odměně**,
 - konečný počet interakcí mezi agentem a prostředím – **mění způsob výpočtu skóre na průměrování** a **požaduje vyvážená prostředí**,
 - složitostní funkce inspirovaná Levinovou *Kt* zachová Occamovu břitvu.
- **Fyzický čas a adaptivní testovací procedura.**

Insa-Cabrera et al. (2011):

- (příliš) zjednodušený *prototyp a demo experimenty*,
- **druhově specifická rozhraní** umožňují použít stejný test.

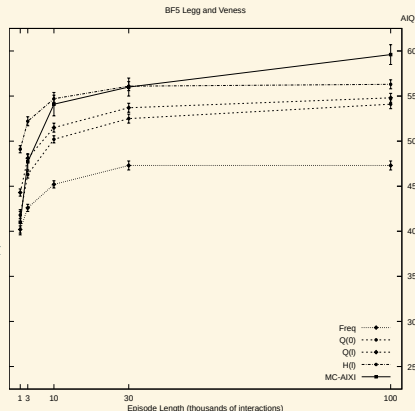
Test algoritmického IQ IV

● Open source implementace:

- odhad pro 200 000 programů prostředí z 10 000 vzorků (pro MC-AIXI 1 000 vzorků),
- test po 1, 3, 10, 30 a 100 tisících interakcích,
- výchozí nastavení BF stroje.

● Agenti v testu:

- *random* – random behavior,
- *freq* – chooses action with highest average reward,
- Q_0 – basic Q-learning,
- $Q\lambda$ – Q-learning with eligibility traces,
- $HLQ\lambda$ – Q-learning with automatic learning rate,
- *MC-AIXI* – wrapper for Monte Carlo approximation of AIXI.



Obrázek: AIQ různých agentů

Diskuze

- Skupiny přístupů vyhodnocující inteligenci s odlišnými předpoklady:
 - ① **Úspěch ve složité úloze je postačující podmínkou pro inteligenci** (Turingův test) *postačitelnost lze zpochybnit pro nelidské subjekty.*
 - ② **Explicitní ověření úspěchu v jednoduchých i složitých prostředích je nutné pro přiznání inteligence** (AIQ test) *dostatečně robustní pro lidi, umělé systémy a zvířata.*

Zdroje I

ANDERSON, J. R. et al. An Integrated Theory of the Mind. *Psychological Review*. 2004, vol. 111, no. 4, s. 1036--1060. ISSN 0033-295X. doi: 10.1037/0033-295X.111.4.1036.

BESOLD, T. R. – HERNÁNDEZ-ORALLO, J. – SCHMID, U. Can Machine Intelligence be Measured in the Same Way as Human intelligence? *KI – Künstliche Intelligenz*. 2015, vol. 29, no. 3, s. 291--297. doi: 10.1007/s13218-015-0361-4.

BRINGSJORD, S. – SCHIMANSKI, B. What Is Artificial Intelligence? Psychometric AI as an Answer. In GOTTLOB, G. (Ed.) *Proceedings of the 18th International Joint Conference on Artificial Intelligence (IJCAI'03)*, s. 887--893, 2003.

Zdroje II

- BURGE, T. Individualism and the Mental. *Midwest studies in philosophy*. 1979, vol. 4, no. 1, s. 73--121. ISSN 1475-4975. doi: 10.1111/j.1475-4975.1979.tb00374.x.
- MEY, M. *The Cognitive Paradigm*. Chicago and London: University of Chicago Press, 1992. doi: 10.1007/978-94-009-7956-7. ISBN 0-226-14259-0.
- DESCARTES, R. *Rozprava o metodě*. Praha: Svoboda, 3. vyd., 1637. 1992 (poprvé publikováno 1637).
- GOERTZEL, B. Artificial General Intelligence: Concept, State of the Art, and Future Prospects. *Journal of Artificial General Intelligence*. 2014, vol. 5, no. 1, s. 1--48. doi: 10.2478/jagi-2014-0001.

Zdroje III

GOERTZEL, B. Toward a Formal Characterization of Real-World General Intelligence. In BAUM, E. – HUTTER, M. – KITZELMANN, E. (Ed.) *Proceedings of the 3rd International Conference on Artificial General Intelligence (AGI 2010), Lugano, Switzerland*, 11 / *Advances in Intelligent Systems Research*, s. 19--24, Amsterdam-Beijing-Paris, 2010. Atlantis Press. doi: 10.2991/agi.2010.17. ISBN 978-90-78677-36-9.

HARNAD, S. Other Bodies, Other Minds: A Machine Incarnation of an old Philosophical Problem. *Minds and Machines*. 1991, vol. 1, no. 1, s. 43--54. ISSN 0924-6495. doi: 10.1007/BF00360578.

Zdroje IV

- HAVEL, I. M. Přirozené a umělé myšlení jako filozofický problém. In MAŘÍK, V. – ŠTĚPÁNKOVÁ, O. – LAŽANSKÝ, J. (Ed.) *Umělá inteligence 3*. Praha: Akademia, 1. vyd., 2001. s. 17--75. ISBN 80-200-0472-6.
- HERNÁNDEZ-ORALLO, J. *The Measure of All Minds*. Cambridge: Cambridge University Press, 1. vyd., 2017. doi: 10.1017/9781316594179. ISBN 978-1-10715-301-1.
- HERNANDEZ-ORALLO, J. Beyond the Turing Test. *Journal of Logic, Language and Information*. 2000, vol. 9, no. 4, s. 447--466. ISSN 0925-8531. doi: 10.1023/A:1008367325700.
- HERNÁNDEZ-ORALLO, J. On environment difficulty and discriminating power. *Autonomous Agents and Multi-Agent Systems*. 2015, vol. 29, no. 3, s. 402--454. ISSN 1387-2532. doi: 10.1007/s10458-014-9257-1.

Zdroje V

HERNÁNDEZ-ORALLO, J. – DOWE, D. L. Measuring Universal Intelligence: Towards an Anytime Intelligence Test. *Artificial Intelligence*. 2010, vol. 174, no. 18, s. 1508--1539. ISSN 0004-3702. doi: 10.1016/j.artint.2010.09.006.

HIBBARD, B. Bias and No Free Lunch in Formal Measures of Intelligence. *Journal of Artificial General Intelligence*. 2009, vol. 1, no. 1, s. 54--61. doi: 10.2478/v10229-011-0004-6.

HUTTER, M. One Decade of Universal Artificial Intelligence. In WANG, P. – GOERTZEL, B. (Ed.) *Theoretical Foundations of Artificial General Intelligence, 4 / Atlantis Thinking Machines*. Paris: Atlantis Press, 2012. s. 67--88. doi: 10.2991/978-94-91216-62-6_5. ISBN 978-94-6239055-3.

Zdroje VI

HYSLOP, A. Other Minds. In ZALTA, E. N. (Ed.) *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, spring 2014. vyd., 2014. Dostupné z:
<<http://plato.stanford.edu/archives/spr2014/entries/other-minds/>>.

INSA-CABRERA, J. et al. Comparing Humans and AI Agents. In SCHMIDHUBER, J. – THÓRISSON, K. R. – LOOKS, M. (Ed.) *Proceedings of the 4th International Conference on Artificial General Intelligence (AGI 2011), Mountain View, USA, 6830 / Lecture Notes in Artificial Intelligence*, s. 122--132, Berlin, 2011. Springer. doi: 10.1007/978-3-642-22887-2_13. ISBN 978-3-642-22887-2.

Zdroje VII

- KRIPKE, S. A. *Naming and necessity*. Cambridge: Harvard University Press, 1972. ISBN 978-0-674-59846-1.
- LEGG, S. – HUTTER, M. Universal Intelligence: A Definition of Machine Intelligence. *Minds and Machines*. Dec 2007, vol. 17, no. 4, s. 391--444. ISSN 0924-6495. doi: 10.1007/s11023-007-9079-x.
- LEGG, S. – VENESS, J. An Approximation of the Universal Intelligence Measure. In DOWE, D. L. (Ed.) *Algorithmic Probability and Friends. Bayesian Prediction and Artificial Intelligence, 7070 / Lecture Notes in Computer Science*. Berlin, Heidelberg: Springer, 2013. s. 236--249. doi: 10.1007/978-3-642-44958-1_18.

Zdroje VIII

- MÜLLER, U. *dev/lang/brainfuck-2.lha in Aminet* [online]. 9. 6. 1993. [cit. 2017-06-26]. Dostupné z: <<http://aminet.net/package.php?package=dev/lang/brainfuck-2.lha>>.
- PSTRUŽINA, K. Turing v Searleově čínské místnosti. *Acta oeconomica Pragensia*. 2003, vol. 11, no. 8, s. 9--16.
- PUTNAM, H. *The Meaning of 'Meaning'*, s. 215--271. Cambridge: Cambridge University Press, 1975. ISBN 978-0-521-20668-6.
- SCHWEIZER, P. The Externalist Foundations of a Truly Total Turing Test. *Minds and Machines*. Aug 2012, vol. 22, no. 3, s. 191--212. ISSN 0924-6495. doi: 10.1007/s11023-012-9272-4. <http://www.research.ed.ac.uk/portal/en/persons/paul-schweizer%2894664317-391f-4731-9125-577d61b109ba%29.html>.

Zdroje IX

- SEARLE, J. R. Minds, Brains, and Programs. *Behavioral and Brain Sciences*. 1980, no. 3, s. 417--457. doi: 10.1017/S0140525X00005756.
- SUN, R. The Importance of Cognitive Architectures: An Analysis Based on CLARION. *Journal of Experimental & Theoretical Artificial Intelligence*. 2007, vol. 19, no. 2, s. 159--193. ISSN 0952-813X. doi: 10.1080/09528130701191560.
- THOMSEN, K. The Cerebellum in the Ouroboros Model, the “Interpolator Hypothesis”. In SHIMIZU, S. – BOSSOMAIER, T. (Ed.) *Proceedings of the 5th International Conference on Advanced Cognitive Technologies and Applications (COGNITIVE 2013), Valencia, Spain*, s. 37--41, Wilmington, 2013. IARIA. Dostupné z: <<http://www.thinkmind.org/download.php>>

Zdroje X

articleid=cognitive_2013_2_30_40069>. ISBN 978-1-61208-273-8.

THÓRISSON, K. R. et al. Towards Flexible Task Environments for Comprehensive Evaluation of Artificial Intelligent Systems and Automatic Learners. In BIEGER, J. – GOERTZEL, B. – POTAPOV, A. (Ed.) *Proceedings of the 8th International Conference on Artificial General Intelligence (AGI 2015), Berlin, Germany, 9205 / Lecture Notes in Artificial Intelligence*, s. 187--196, Berlin, 2015. Springer. doi: 10.1007/978-3-319-21365-1_20. ISBN 978-3-319-21364-4.

Zdroje XI

- THÓRISSON, K. R. et al. Why Artificial Intelligence Needs a Task Theory And What It Might Look Like. In STEUNEBRINK, B. – WANG, P. – GOERTZEL, B. (Ed.) *Proceedings of the 9th International Conference on Artificial General Intelligence (AGI 2016), New York, USA*, 9782 / *Lecture Notes in Artificial Intelligence*, s. 118--128, New York, 2016. Springer. doi: 10.1007/978-3-319-41649-612. ISBN 978-3-319-41649-6.
- TURING, A. M. Computing Machinery and Intelligence. *Mind*. 1950, vol. 59, no. 236, s. 433--460. ISSN 0026-4423.
- VADINSKÝ, O. Přehled přístupů k vyhodnocování inteligence umělých systémů. *Acta Informatica Pragensia*. 2018, vol. 7, no. 1, s. 74--103. ISSN 1805-4951. doi: 10.18267/j.aip.115.

Zdroje XII

WANG, P. On Defining Artificial Intelligence. *Journal of Artificial General Intelligence*. 2019, vol. 2, no. 10, s. 1--37. doi: 10.2478/jagi-2019-0002.