

# Sztuczna inteligencja. Wykład wstępny

Paweł Rychlikowski

Instytut Informatyki UWr

3 marca 2025

- Zapraszam w każdy czwartek!
- Mile widziane przerywanie pytaniami, itd.

- Mamy ćwiczenia i pracownie, standardowo w rytmie:  
 $5 * (P + P + C)$
- Na każdą parę pracowni będzie lista pracowniowa, na każde ćwiczenia będzie lista ćwiczeniowa
- Dodatkowo można pisać nieobowiązkowy projekt (z grubsza w czasie P5)
- Zwyczajowo uczestnictwo **Codig Game Spring/Summer Challenge** jest również częścią naszego przedmiotu

Treścią pracowni są indywidualne konsultacje:

- dotyczące zadań, które Student(ka) wykonał(a)
- **oraz tych, z którymi ma kłopot**

Kilka drobiazgów:

- Sprawdzanie może być selektywne – prowadzący nie musi oglądać wszystkich zadań (ale za wszystkie 'zadeklarowane' powinien przyznać punkty)
- Student ma prawo zadawać różne pytania, prosić o wskazówki, o pomoc przy debugowaniu. itp
- Prowadzący decyduje, jak wygląda opowiadanie o zadaniu. Może na przykład poprosić o zreferowanie rozwiązania, albo zadać pytanie o jakiś fragment, albo przeczytać kod i stwierdzić, że wszystko jasne, albo ...

## Prośba od prowadzących

Opowiadając o zadaniu nie trzeba mówić o wszystkich funkcjach. Znaczną część rutynowych rzeczy można (należy) pomijać

- **Każde zadanie powinno mieć b. krótki opis (w pierwszym komentarzu)**
- W opisie należy umieścić najważniejsze informacje o zadaniu, przykładowo (dla zadania szachowego):

W rozwiązaniu używam standardowego algorytmu BFS. Sytuację na planszy pamiętam jako string, postaci `''biały król jest na [pole1], czarny na [pole2], a wieża na [pole3]''`. Przy teście 3 konieczne jest 32 GB pamięci (inaczej mamy segfault)

# O sprawdzaniu zadań

## Sprawdzaczka

```
python validator.py zad1 solution.exe
```

```
python validator.py zad1 python zad1_sol.py
```

Możliwe inne rozwiązania (np. wykorzystanie Coding Game)

- Jest językiem sprawdzaczki (więc trzeba go mieć zainstalowanego)
- Nadaje się do większości zadań
- Może warto wykorzystać te zajęcia, żeby się go poduczyć
- Ma wiele bibliotek wspierających różne zadania związane z AI

**Ale ogólnie nie wymuszamy żadnego języka programowania**



# Jak nie Python, to co?

Inne języki często używane w kontekście AI:

- **R**: wsparcie dla ML (zły wybór na naszym przedmiocie)
- **C++, Java, Go, Rust, ...**: szybkość symulacji, przeszukiwania, efektywne zarządzanie pamięcią
- **Lisp, Prolog**: dobre dopasowanie do **niektórych** zadań AI

- Zasadniczo zakładamy, że rozwiązujemy ją bazując na wstępnej wiedzy
- ale oczywiście można spytać o dowolne zadanie, również dzisiaj po wykładzie

- System deklaryacyjny (przed zajęciami mówimy, co umiemy i za te zadania dostajemy punkty)
- Dodatkowa zachęta dla rozwiązujących: bonus +2 za 50%

- Za prezentację (i wcześniejsze napisanie) programów studenci dostają punkty
- Podobnie za deklarację zadań na ćwiczeniach

## **Wszystkie punkty są równe!**

Sumujemy je i wyznaczamy ocenę (szczegóły w regulaminie)

- Stuart Russel, Peter Norvig, Artificial Intelligence. A Modern Approach. 3rd edition (w Internecie leży pdf)
- Fajna, ale 1100 stron.
- **Jest już czwarte wydanie!** (jak ktoś chce kupić, to koniecznie w miękkiej okładce!)
- Opcja *istotnie* tańsza (po polsku): AIMA podzielona na dwa tomy, wydana przez Helion



# Sztuczna inteligencja. Definicja

## Definicja (krótka)

Zdolność komputera (programu) do wykonywania zadań powszechnie kojarzonych z zachowaniem inteligentnym (ludzi lub zwierząt).

## Definicja długa

Spis treści wybranego podręcznika o Sztucznej inteligencji

# Sztuczna inteligencja (definicja długa)

## Definicja długa – spis treści

Rozwiązywanie problemów przez przeszukiwanie, rozwiązywanie więzów, przeszukiwanie z oponentem, logika w opisie świata, wnioskowanie w logice, planowanie, modelowanie niepewności, reprezentacja wiedzy, wnioskowanie przy niepewności, podejmowanie decyzji, uczenie się z przykładów, uczenie się ze wzmocnieniem, przetwarzanie języka naturalnego, rozpoznawanie wzorców w obrazach i dźwiękach, robotyka, procedural content generation.

## Przedmioty:

1. Machine Learning
2. Neural Networks: Theory and Practice
3. Modele językowe
4. Advanced Data Mining
5. AI for games (i inne z tej okolicy)



# Definicja na bazie wiedzy „portalowej”

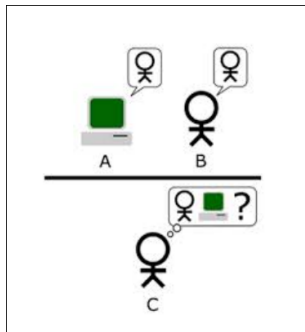
Ostatnio fraza **Sztuczna inteligencja** bywa używana w specyficznych kontekstach.

Przykłady:

- Sztuczna inteligencja zabierze pracę twórcom?
- Sztuczna inteligencja stworzyła zdjęcia z imprezy...
- Zapytali sztuczną inteligencję o górnictwo w Polsce. Co im odpowiedziała...

Bardzo często fraza **Sztuczna inteligencja** dotyczy tzw. **modeli generatywnych** (czyli tworzących treści cyfrowe, teksty, obrazy, wideo)

# Test Turinga



- Pytanie: Czy maszyny potrafią myśleć?
- Odpowiedź: Możemy to sprawdzić w rozmowie.

## Pytanie

Czy **naprawdę** chcemy, żeby komputer był nierozróżnialny w rozmowie od człowieka?



źródło: tvtropes.org

**Nieludzkie, pożądane cechy:** wiedza, zdolności rachunkowe, modelowanie świata, rozwiązywanie zadań fizycznych, ...

# Red Flag for AI



# Czy komputery mogą przejść test Turinga

Czy umiemy odpowiedzieć na to pytanie w lutym 2025?

# Przykładowy agent konwersacyjny

```
x = input('Ty: ')
print ('Bot:', 'Cześć, miło mi Cię spotkać')
x = input('Ty: ')
print ('Bot:', 'Trudno mi o tym mówić, bowiem mój
          chomik... On chyba ma depresję!')
x= input('Ty: ')
print ('Bot:', 'Dzięki za wsparcie.')
```

# Dialog 1

**[User:]** Cześć, jestem Ania

**[Bot:]** Cześć, miło mi Cię spotkać

**[User:]** Jak dojadę do dworca z pl. Grunwaldzkiego?

**[Bot:]** Trudno mi o tym mówić, bowiem mój chomik... On chyba ma depresję!

**[User:]** Przykro mi...

**[Bot:]** Dzięki za wsparcie.

## Dialog 2

**[User:]** Hej, tu Piotrek

**[Bot:]** Cześć, miło mi Cię spotkać

**[User:]** Gdzie najlepiej studiować informatykę we Wrocławiu?

**[Bot:]** Trudno mi o tym mówić, bowiem mój chomik... On chyba ma depresję!

**[User:]** Kurcze, nie przyszedłem tu rozmawiać o jakichś chomikach!

**[Bot:]** Dzięki za wsparcie.

Zwróćmy uwagę na sarkazm ostatniej wypowiedzi (kropka na końcu!)



# Test Turinga w rzeczywistości

Programy, które **czasami** przechodzą test Turinga, są z nami od bardzo dawna. Istotne przykłady:

- **ELIZA** (1965) – udawała psychoterapeutę (i potrafiła zmylić pacjentów)
- **PARRY** (1972) – udawał osobę z zaburzeniami (i potrafił zmylić psychiatrów)
- **Loebner Prize** (1990 – 2019), w 2008 zwycięzca zmylił 3 z 12 sędziów
- **LaMDA**, (2022) – Blake Lemoine (inżynier w Google) twierdził na podstawie rozmów, że ten model językowy jest świadomy

Welcome to

EEEEEE	LL	IIII	ZZZZZZ	AAAAA
EE	LL	II	ZZ	AA AA
EEEEE	LL	II	ZZZ	AAAAAAA
EE	LL	II	ZZ	AA AA
EEEEEE	LLLLLL	IIII	ZZZZZZ	AA AA

Eliza is a mock Rogerian psychotherapist.  
The original program was described by Joseph Weizenbaum in 1966.  
This implementation by Norbert Landsteiner 2005.

ELIZA: Is something troubling you ?  
YOU: Men are all alike.  
ELIZA: What is the connection, do you suppose ?  
YOU: They're always bugging us about something or other.  
ELIZA: Can you think of a specific example ?  
YOU: Well, my boyfriend made me come here.  
ELIZA: Is it important to you that your boyfriend made you come here ?  
YOU: He says I'm depressed much of the time.  
ELIZA: I am sorry to hear that you are depressed.  
YOU: It's true. I am unhappy.  
ELIZA: Can you explain what made you unhappy ?  
YOU:

## Definicja

**Modelem językowym** nazywamy rozkład prawdopodobieństwa na ciągach słów:

$$p(w_{n+1} | w_1 w_2 \dots w_n)$$

Model językowy może być użyty do generowania tekstu:

- Wybieramy (najbardziej?) prawdopodobne  $w_{n+1}$  i powtarzamy proces

## Przykładowe metody wyboru

- Wybieramy najbardziej prawdopodobny (nudne...)
- Losujemy z  $K$  najbardziej prawdopodobnych następników
- Losujemy wśród najbardziej prawdopodobnych następników, które łącznie mają prawdopodobieństwo  $p$

Modelami językowymi są:

- GPT-1, GPT-2, GPT-3 (po prostu modelują język)
- ChatGPT (GPT 3.5), GPT-4, LaMDA, ... (dodatkowo dotrenowywane do dialogów)
- Nie tylko język:
  - AlphaFold (modeluje białka) – Nobel 2024 z chemii,
  - część AlphaStar (modeluje rozgrywkę StarCrafta),
  - T5 (robi różne rzeczy, w tym tłumaczenie czy klasyfikacja tekstów),
  - GATO (modeluje wiele rzeczy, w tym rozwiązywanie łamigłówek, czy gry na stare konsole)

- Wielkości modeli:
  - GPT-2: 1.5 G-float
  - GPT-3: 175 G-float (podstawa ChatGPT)
  - GPT-4 (?): 100 T-float (w zapowiedziach, raczej dużo mniej)

Czy ilość przechodzi w jakość? (otwarte pytanie)

**Suppose I have 4 apples, and I gave two of them to my sister**, and she gave one of them to her sister, and it started to grow. Would I be able to see how much time it had taken for the apple to grow?

We will first give you a number which you will later use to work out how much time it has taken for the apple to grow:

(Time it took to grow 2 apples, shown as a clock)  $\times$  (Size of apple)  
= "time taken"

What this tells us is that:

$2 \times 3 \times 2 = 44$  (seconds in 24 hours)

So the question is, what does 44 seconds look like?

Let's do a quick experimen

# ChatGPT o naszym przedmiocie

jest tworzenie systemów, które będą potrafiły w sposób autonomiczny podejmować decyzje i wykonywać zadania, które obecnie wymagają udziału człowieka.

P

Dziękuję, ale prosiłbym o krótszą definicję, najlepiej w formie przyśpiewki weselnej



Oczywiście, z przyjemnością pomogę! Oto przyśpiewka weselna na temat sztucznej inteligencji:

Sztuczna inteligencja to magia,  
Komputery, co same się uczą,  
Algorytmy, sieci neuronowe,  
Decyzje podejmują jak człowiek prawdziwy.

Dane to klucz, nauczanie maszynowe,  
Rozpoznawanie obrazów, mowy,  
Sztuczna inteligencja na każdym kroku,  
Świat zmienia, zaskakuje, niesamowity!

# Słowa klucze dla naszego wykładu (2024)

- Przeszukiwanie (problem solving)
- Uczenie się (learning)
- Wnioskowanie (inference, theorem proving)
- Modelowanie wiedzy o świecie

## Główny problem SI

Głównym, nierozwiązanym do tej pory problemem SI jest satysfakcjonujące połączenie tych światów, zwłaszcza wnioskowania (twarda logika), z rozmytą wiedzą pochodzącą z uczenia się.



# Słowa klucze dla naszego wykładu (2025)

- Przeszukiwanie (problem solving)
- Uczenie się (learning)
- Wnioskowanie (inference, theorem proving)
- Modelowanie wiedzy o świecie

## Główny problem SI

Głównym, nierozwiązanym (?) do tej pory problemem SI jest satysfakcjonujące połączenie tych światów, zwłaszcza wnioskowania (twarda logika), z rozmytą wiedzą pochodzącą z uczenia się.

# Łączenie dwóch światów

- Model językowy (jak jest odpowiednio duży) może produkować sensowniejsze rozumowania.
- Można go też do tego zachęcić, dodając do zapytania coś w stylu: **Let us think step by step**
- Model językowy może się również uczyć coraz mądrzejszych rozumowań (OpenAI o1, DeepSeek R1)
- Model językowy może współpracować z zaawansowanymi algorytmami przeszukiwania i formalnymi systemami dowodzącymi (AlphaProof, ten od medalu w IMO)

Gdyby tylko LLM-y były 1000 razy mniejsze i 1000 razy szybsze...

# Dlaczego problemy SI są trudne

Istnieją dwa główne powody

- ① Trudność algorytmiczna:
  - **NP-zupełne** - łamigłówki,
  - **PSPACE-zupełne** - gry,
  - **EXPTIME i wyżej** – systemy wnioskowania
- ② Czasem trudno sformalizować precyzyjnie problem!

# Co się udało Sztucznej inteligencji?

## Gdzie komputery są od nas lepsze?

1. Oczywiście my też wykonujemy algorytmy (np. mnożenie liczb). Tu od zawsze **one** mają przewagę.
  - chyba, że algorytm wykonuje LLM (model językowy), np. praca: GPT Can Solve Mathematical Problems Without a Calculator
2. Komputery radzą sobie w grach: Backgammon (tryktrak) (1992), Szachy (1996), Go (2016), Poker (2017), StarCraft (AlphaStar, 2019?), Stratego (2022), Dyplomacja (2022)
3. Rozpoznawanie mowy (Microsoft/IBM, Switchboard Corpora, około 2016)

# Co się udało Sztucznej inteligencji? (cd)

## Gdzie komputery są od nas lepsze?

1. Rozpoznawanie prostych obrazów (ImageNet, 1mln obrazków, 1000 klas):
  - 2010: około 28% błędów
  - 2015 (4.94% na ImageNet, człowiek: 5.1%)
  - W 2020 było 1.3%, w 2021 było 1.2%, 2023: **0.98%**.
2. Tłumaczenie maszynowe
3. Wygrywanie teleturniejów wiedzy (Watson, Jeopardy, 2011),
4. Rozwiązywanie zadań z Międzynarodowej Olimpiady Matematycznej, **AlphaProof**, okolice srebrnego medalu w 2024

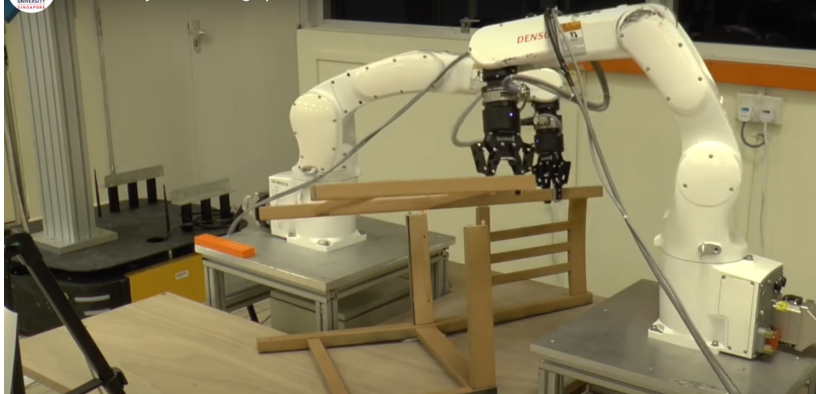
# Co działa niekoniecznie idealnie?

- ❶ Rozpoznawanie mowy (50% błąd na nagraniach *przy koltlecie*)
- ❷ Roboty umiejące nalewać wodę, otwierać drzwi, itd w **nieznanym** środowisku
  - DARPA challenge fails
  - RoboCup (zadanie domowe: Finale RoboCup 2023 League Kid Size)
- ❸ Gra w Brydża (choć w 2022 NukAI program komputerowy wygrał „turniej brydżowy”, ale tylko jako rozgrywający i tylko 3NT)
- ❹ „Ludzka” rozmowa na dowolny temat (??)

# The IKEA test



Robot by NTU Singapore builds an IKEA chair



(film z roku 2018)

## Zadanie (na przyszłość)

Para robotów ma złożyć **dowolny** mebel IKEA, korzystając jedynie z instrukcji umieszczonej w pudełku.

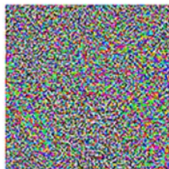
# Ataki na rozpoznawanie obrazów



"panda"

57.7% confidence

+  $\epsilon$



=



"gibbon"

99.3% confidence

An image of a panda, when combined with an adversarial input, can convince a classifier that it's looking at a gibbon. IMAGE: OPENAI

Źródło: Slight Street Sign Modifications Can Completely Fool Machine Learning Algorithms



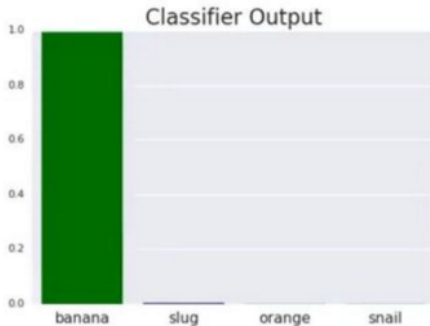
Sztucznie wygenerowany obraz, maksymalizujący **tosterowość**.



# Co możemy zrobić z tym obrazkiem?

- Możemy go pokazywać sieci.
- Ale wklejając go analogowo, nie cyfrowo.
- Zobacz pracę: Adversarial Patch, T. Brown i inni, 2017

# Tostery i banany

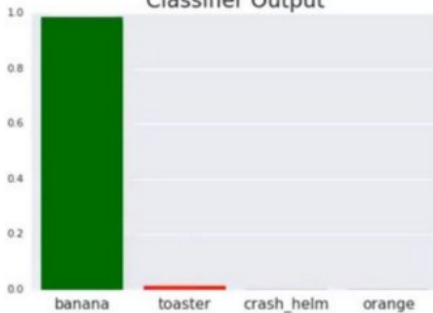


# Tostery i banany

Classifier Input



Classifier Output

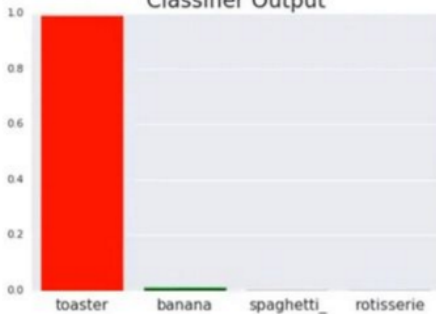


# Tostery i banany

Classifier Input



Classifier Output



Ograniczenia prędkości do **45 mil na godzinę**



Figure 1: The left image shows real graffiti on a Stop sign, something that most humans would not think is suspicious. The right image shows our a physical perturbation applied to a Stop sign. We design our perturbations to mimic graffiti, and thus “hide in the human psyche.”

Źródło: Robust Physical-World Attacks on Deep Learning Visual Classification

Najpierw zajmiemy się **przeszukiwaniem**, które jest jednym z podstawowych narzędzi AI.

# Problem solving by searching. Intuicje

## Przykład 1. Wyznaczanie trasy



## Przykład 2. Wyznaczanie sekwencji działań

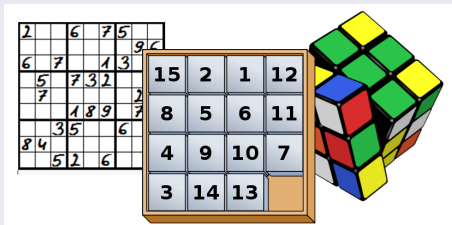
Kohler (1945): monkey and banana problem.



Kohler observed that chimpanzees appeared to have an insight into the problem before solving it



## Przykład 3. Rozwiązywanie łamigłówek



## Uwaga

**Problemy zabawkowe** (toy problems) są częstym narzędziem w AI.