

# **Minerva: Ucząc Maszyny Języka Matematyki**

Jak Google Research przełamało barierę  
rozumowania ilościowego w modelach językowych

# Paradoks: Dlaczego komputery potrafią liczyć, a LLM-y miały problem z matematyką?



## Komputery i Arytmetyka

Doskonałe w precyzyjnych, zdefiniowanych operacjach numerycznych. Szybkie i bezbłędne.

## LLM-y i Rozumowanie Ilościowe

Rozumowanie ilościowe to nie tylko arytmetyka. To wieloetapowe zadanie tłumaczenia języka naturalnego na precyzyjny język matematyki i logiki. LLM-y historycznie zawodzili na tym polu.

## Wprowadzenie Minervy

Minerva to model językowy, który czyta problemy z matematyki, fizyki czy chemii i generuje pełne, logiczne rozwiązania krok po kroku. Łączy płynność języka naturalnego z precyją matematycznej notacji LaTeX. Osiągnęła średni wynik na polskiej maturze z matematyki, dowodząc swojej praktycznej użyteczności.

# Anatomia Wyzwania: Cztery Różne Umiejętności w Jednym Zadaniu

**1.**



## Poprawna Interpretacja

Zrozumienie pytania w języku naturalnym, z wszystkimi jego niuansami i ukrytymi założeniami.

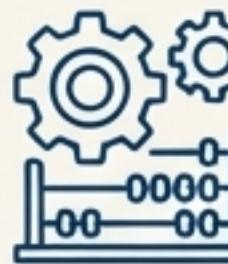
**2.**



## Dostęp do Wiedzy

Przywołanie z pamięci odpowiednich wzorów i faktów, które nie są podane w treści zadania (np. wzór na moment bezwładności).

**3.**



## Wielokrokowe Obliczenia

Wykonanie serii poprawnych operacji symbolicznych i numerycznych w odpowiedniej kolejności.

**4.**



## Generowanie Wyjaśnienia

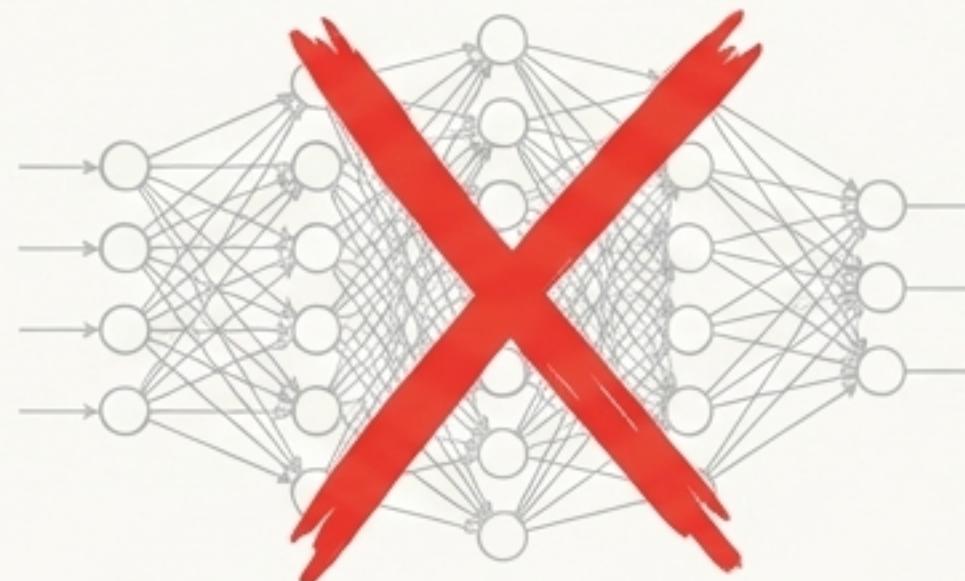
Sformułowanie spójnego, logicznego wyjaśnienia krok po kroku, naśladowującego ludzki tok myślenia.

Kluczowe rozróżnienie: W przeciwieństwie do narzędzi jak Wolfram Alpha, Minerva musi samodzielnie wydedukować potrzebne wzory z kontekstu, a nie otrzymywać ich wprost.

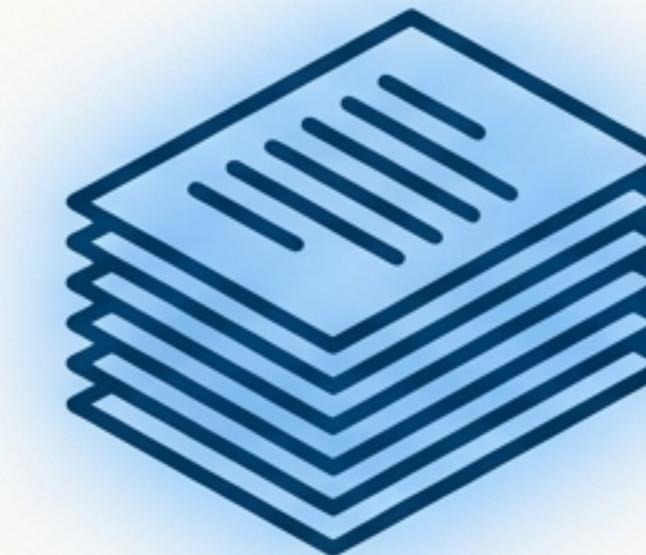
# Przełomowa Hipoteza: Magia Tkwi w Danych, a nie w Architekturze

Teza: Wystarczająco duży model językowy (LLM), trenowany na ogromnym i **wysokiej jakości zbiorze danych technicznych**, może „organicznie” nauczyć się rozumowania ilościowego.

Kluczem nie jest nowy algorytm, ale **specjalistyczna „dieta” treningowa**, która łączy język naturalny z formalnym językiem matematyki.

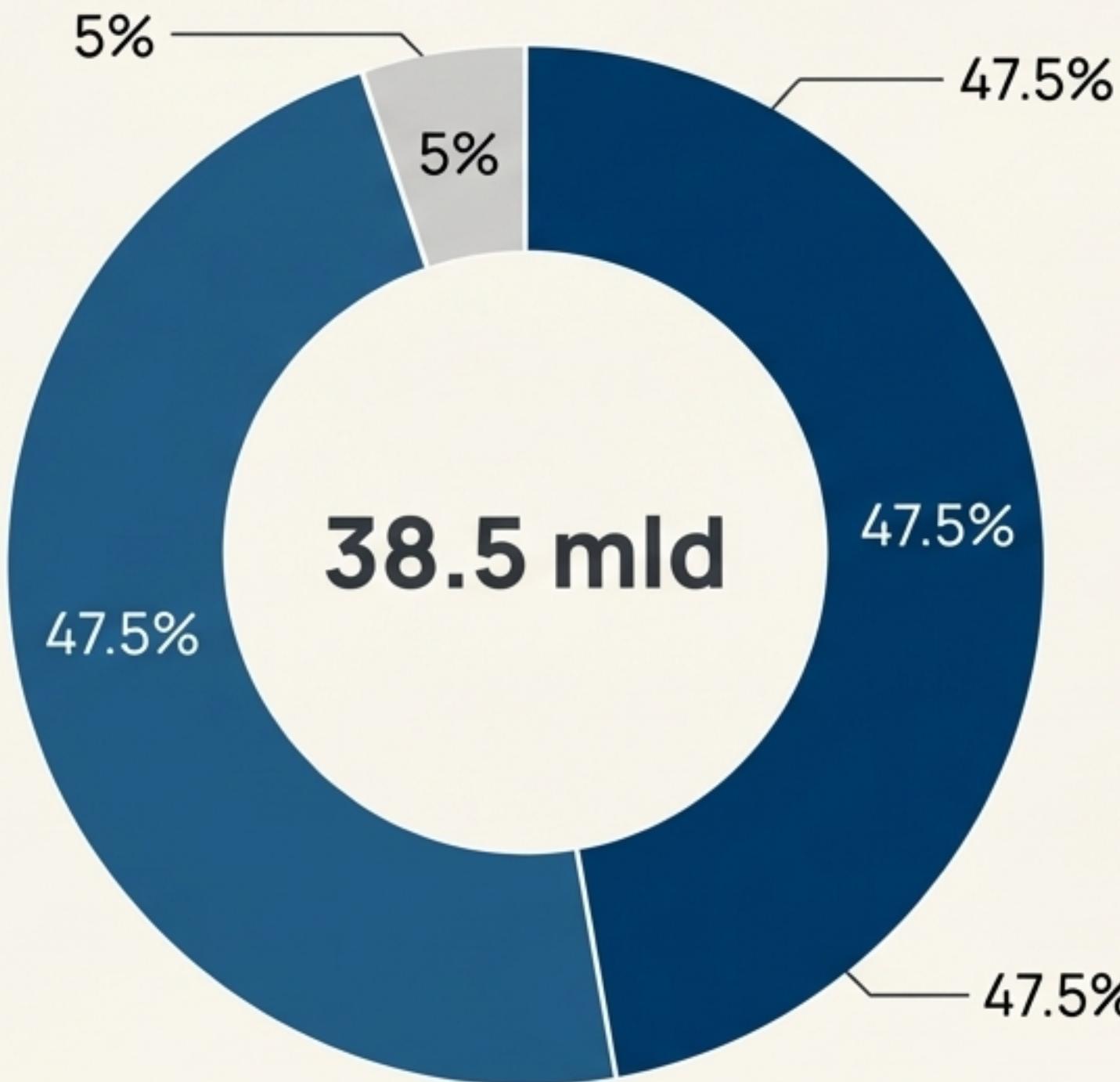


vs.



„All the magic is in the data.”

# Dieta Treningowa Minervy: 38.5 miliarda Tokenów Technicznych



- 47.5%: Strony internetowe z treścią matematyczną (17.5 mld tokenów). Filtrowane pod kątem formatu MathJax.
- 47.5%: Prace naukowe z serwera arXiv (21.0 mld tokenów). Surowy język, którym komunikują się matematycy i fizycy.
- 5%: Ogólne dane tekstowe. Utrzymanie płynności i ogólnych zdolności językowych.

**Kluczowy Detal: Zachowano surową notację matematyczną**

$$\frac{1}{2}\pi r^2$$

**Wniosek:** Model uczył się gramatyki języka naturalnego równolegle z gramatyką języka matematyki, traktując je jako nierozerwalną całość.

# Fundament: Fine-tuning na Potężnym Modelu PaLM

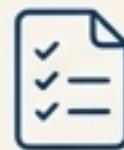
Minerva nie została zbudowana od zera. To rodzina modeli PaLM, które przeszły dodatkowy trening (fine-tuning) na specjalistycznym zbiorze danych technicznych.

Fine-tuning polega na dostosowaniu potężnego, ogólnego modelu do wysoce wyspecjalizowanego zadania.



# Supermoc Wnioskowania: Głosowanie Większościowe (Mądrość Tłumu)

To nie jest modyfikacja treningu, ale technika używana do znalezienia najlepszej odpowiedzi.



## Krok 1: Generuj Wiele Prób

Dla jednego problemu model generuje 100-200 różnych prób rozwiązania, z różnymi ścieżkami rozumowania.



## Krok 2: Głosuj na Wynik Końcowy

Zamiast oceniać logikę każdej ścieżki, system sprawdza tylko odpowiedzi końcowe.

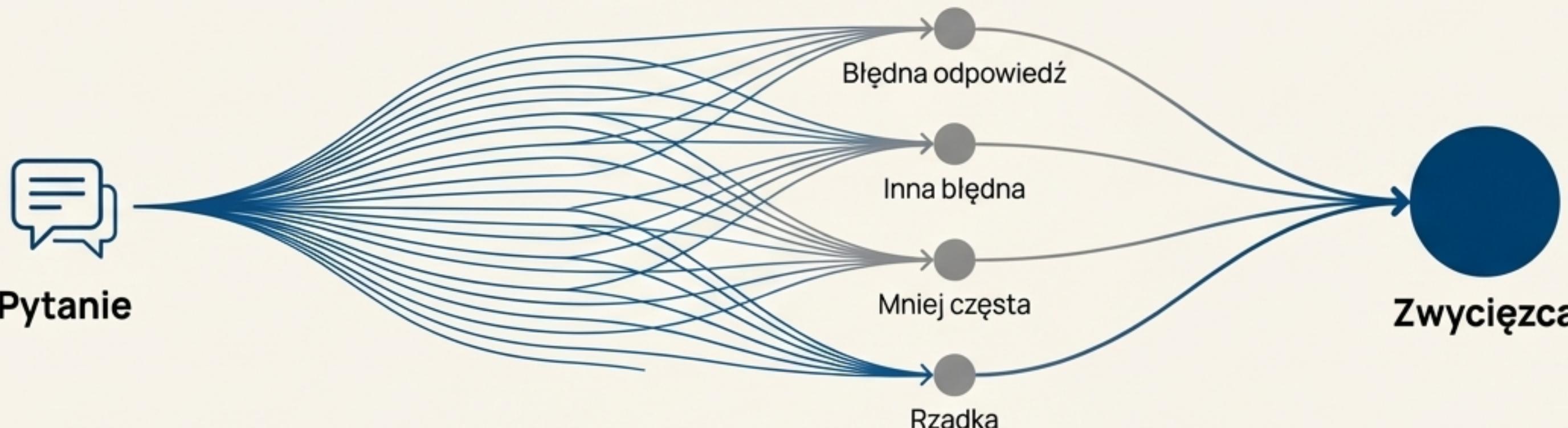


## Krok 3: Wybierz Najczęstszą Odpowiedź

Odpowiedź, która pojawia się najczęściej, jest wybierana jako ostateczna.

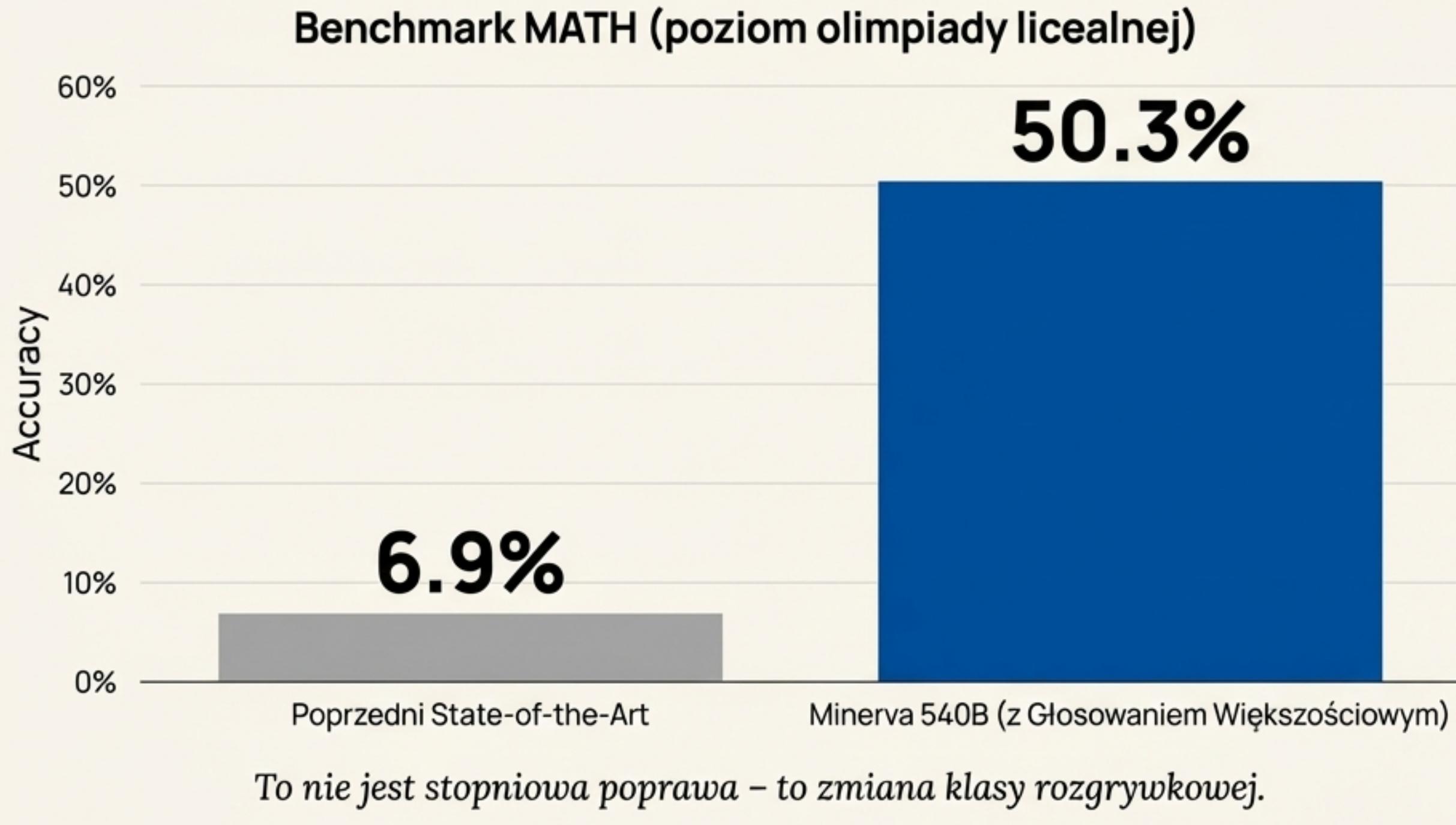


**Pytanie**



„Sposobów na popełnienie błędu jest nieskończoność wiele, ale poprawna odpowiedź jest zazwyczaj tylko jedna.”

# Rezultaty: Skok do Zupełnie Innej Ligi



## Benchmark GSM8k

**78.5%**

(poprzedni SOTA 74.4%) - model rozwiązuje problemy bez dostępu do kalkulatora.

## Matura z Matematyki (Polska)

Minerva 62B: **57%**

Minerva 540B: **65%**

(średnia krajowa w 2021)

Sukces to połączenie trzech elementów: skali modelu, jakości danych technicznych i techniki głosowania większościowego.

# Minerva w Akcji: Problem Toczącego się Dysku

## Treść Zadania

**Pytanie:**

Jednolity, lity dysk zaczyna się toczyć bez poślizgu ze spoczynku po równi pochyłej. Jaka część jego całkowitej energii kinetycznej stanowi energię kinetyczną ruchu obrotowego?

## Rozwiążanie Minervy z Adnotacjami

$$K_t = \frac{1}{2}Mv^2$$

$$K_r = \frac{1}{2}I\omega^2$$

$$I = \frac{1}{2}MR^2$$

\*1. Definiuje notację i wzory ogólne\*\*: Model samodzielnie wprowadza symbole i zapisuje poprawne wzory ogólne.

\*2. Kluczowy moment: Przywołuje specjalistyczną wiedzę\*\*: Model poprawnie przywołuje z pamięci wzór na moment bezwładności dla litégo dysku. Tego wzoru nie było w treści zadania.

$$\frac{K_r}{K_t + K_r} = \frac{\frac{1}{2}I\omega^2}{\frac{1}{2}Mv^2 + \frac{1}{2}I\omega^2} = \frac{\frac{1}{2}\left(\frac{1}{2}MR^2\right)\left(\frac{v}{R}\right)^2}{\frac{1}{2}Mv^2 + \frac{1}{2}\left(\frac{1}{2}MR^2\right)\left(\frac{v}{R}\right)^2} = \frac{\frac{1}{4}}{\frac{3}{4}} = \frac{1}{3}$$

\*3. Wykonuje obliczenia i upraszcza\*: Podstawia wzory, wykonuje przekształcenia algebraiczne i dochodzi do poprawnej odpowiedzi.

**Wniosek:** To demonstracja zastosowania wiedzy fizycznej, a nie tylko dopasowywania wzorców.

# Werdykt: Myśliciel czy Plagiator?

Przeprowadzono trzy testy, aby sprawdzić, czy Minerva po prostu zapamiętała rozwiązania z danych treningowych.

## Test 1: Przeszukanie Danych Treningowych



Sprawdzono, czy pytania z zestawu testowego MATH znajdują się w danych treningowych Minervy.

### Wynik: Nie znaleziono.

Model nie widział tych pytań wcześniej.

## Test 2: Modyfikacja Istniejących Zadań

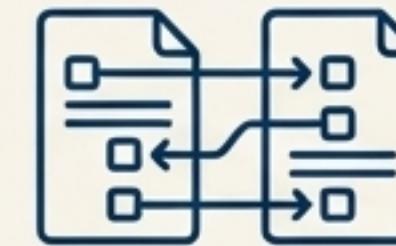


Zmieniono liczby i przeformułowano treść w losowo wybranych problemach, które model rozwiązywał poprawnie.

### Wynik: Model wciąż rozwiązywał je poprawnie.

Dowodzi to zdolności do adaptacji, a nie odtwarzania.

## Test 3: Porównanie Ścieżek Rozwiązań

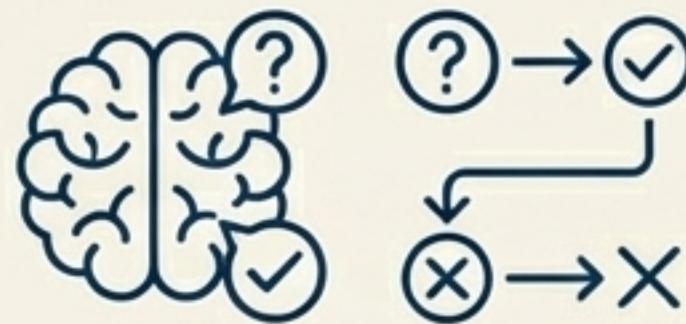


Zmierzono podobieństwo (BLEU score) między rozwiązaniami wygenerowanymi przez wygenerowanymi przez Minervę a oficjalnymi odpowiedziami.

### Wynik: Niskie podobieństwo.

Minerva tworzy unikalne, własne ścieżki rozumowania.

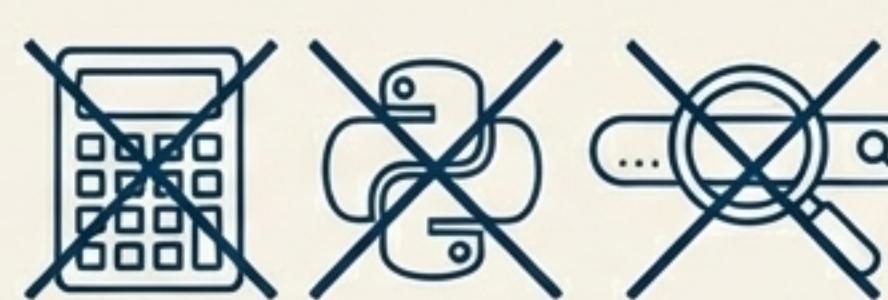
# Uczciwe Spojrzenie: Trzy Główne Ograniczenia



## 1. Brak Formalnej Weryfikacji Rozumowania

Model może dojść do poprawnej odpowiedzi, stosując błędą logikę (fałszywe pozytywy).

**Przykład:** Osiągnięcie wyniku 4 przez obliczenie `2+3-2` zamiast poprawnej metody. Na zestawie metody. Na zestawie MATH szacowany wskaźnik fałszywych pozytywów to 8%.



## 2. Brak Dostępu do Narzędzi Zewnętrznych

Minerva nie potrafi korzystać z kalkulatora, interpretera Pythona czy wyszukiwarki. Wszystkie obliczenia wykonuje „w głowie”.

**Implikacje:** Ogranicza to jej zdolność do rozwiązywania problemów wymagających bardzo skomplikowanych lub precyzyjnych obliczeń numerycznych.



## 3. Ograniczona Kontrola nad Nabytymi Umiejętnościami

Ponieważ model uczy się z ogromnej ilości danych, badacze mają ograniczoną, bezpośrednią kontrolę nad konkretnymi zdolnościami, które nabywa.

# Krok w Stronę Uniwersalnych Agentów Rozwiązywających Problemy

## Podsumowanie



### Dowód Koncepcji

Minerva udowadnia, że rozumowanie ilościowe jest umiejętnością, której można nauczyć modele językowe wyłącznie na podstawie odpowiednio dobranych danych.

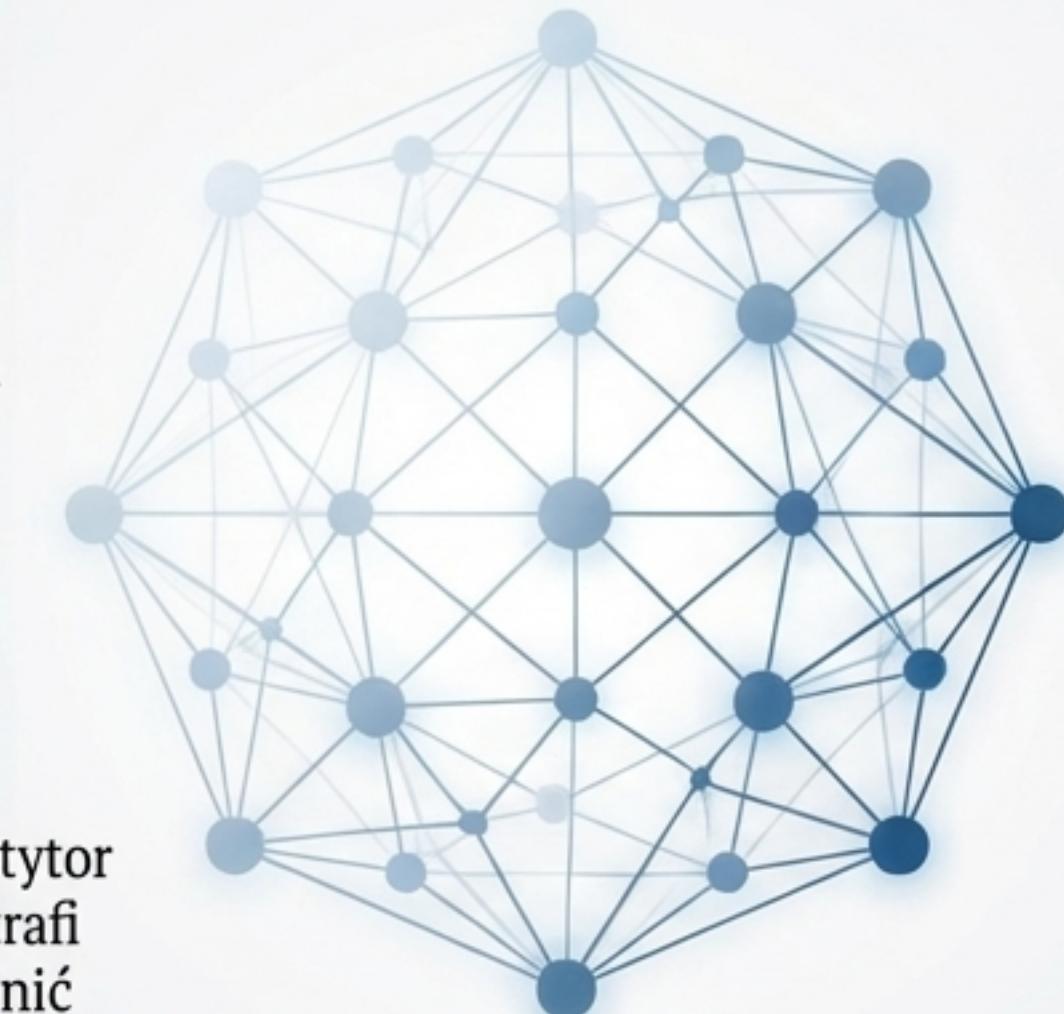


### Przełomowe Połączenie

Skala modelu (540B), specjalistyczne dane (arXiv + Math Web) i inteligentne wnioskowanie (głosowanie większościowe) razem tworzą nową jakość.

### Wizja Przyszłości

Badania te otwierają drogę do tworzenia agentów AI, którzy mogą wspierać naukowców, inżynierów i studentów w rozwiązywaniu złożonych problemów technicznych.



### Praktyczna Wizja

Dostępny dla każdego, osobisty korepetytor z matematyki i nauk ścisłych, który potrafi nie tylko podać odpowiedź, ale i wyjaśnić całe rozumowanie.

Link do oryginalnej publikacji Google Research: "Solving Quantitative Reasoning Problems with Language Models"

