



Metadane w logice artefaktu i w tekście jako sygnał treningowy dla modelu językowego

Co dokładnie znaczy „metadane” w Twoim kontekście

Skoro doprecyzowałeś, że chodzi o **metadane na poziomie logiki artefaktu i tekstu, a nie pliku** (np. nie: „autor PDF”, „data modyfikacji”), to mówimy o takich elementach, które **są zapisane w treści** albo **wynikają z formy i stylu** i które niosą informację „jak czytać, jak oceniać, jak się zachować” — obok informacji „o czym jest świat”.

W Twoich załączonych artefaktach (trzy dokumenty merytoryczne + jeden dokument meta-analityczny o wartości dodanej) widać bardzo wyraźnie, że tekst jest projektowany jako **nośnik procedur epistemicznych** (definicja → hipotezy → falsyfikacja → krytyka) oraz jako **interfejs do dalszego przetwarzania** (atomizacja, rubryki, warianty kontrastowe).

Żeby odpowiedzieć na Twoje „jakie metadane”, najwygodniej potraktować je jako **warstwy sygnału uczącego**:

- **metadane jawne**: dosłownie występujące etykiety/formaty (np. H1–H5, „Definicja robocza (falsyfikowalna)”, tabela kryteriów dowodu),
- **metadane strukturalne**: organizacja argumentu i segmentacja (nagłówki, „sekcja metodologii”, „jak to sfalsyfikować”, kontrasty),
- **metadane epistemiczne**: informacja o statusie twierdzeń (hipoteza/założenie/wniosek/ograniczenie),
- **metadane pragmatyczne**: ton, ostrożność, hedging (np. „interpretacyjnie kruche”, „to nie dowód”),
- **metadane latentne**: powtarzalne wzorce w dystrybucji tokenów (styl recenzencyjny, rytm list i kwalifikatorów), które model uczy się jako korelacje.

Tak rozumiane metadane są bardzo bliskie temu, co w praktykach **instruction tuning / alignment** traktuje się jako „sygnał zachowania”, a nie „wiedzę o świecie”. 1

Konkretnie metadane tekstowe w Twoich plikach: przykłady „co jest czym”

Poniżej wyodrębniam najważniejsze klasy metadanych *na Twoim materiale* (czyli nie ogólnie), wraz z tym, dlaczego to jest ważne treningowo.

Metadane jawne: etykiety, rubryki, definicje testowalności

1) Rubryka „dowód formalny vs empiryczny”

W dokumencie oceniającym formalność „dowodu” masz wprost zdefiniowane, co znaczy dowód formalny (język teorii, aksjomaty, reguły inferencji, teza) i jest zrobiona tabela zgodności kryteriów. To są metadane typu: *standard oceny + checklista*.

Treningowo taki element jest niemal gotową „rubryką scoringową” do uczenia modelu roli recenzenta.

2) Definicje robocze i hipotezy oznaczone jako falsyfikowalne

W „Paradoksie Marii” definicja jest explicite „falsyfikowalna” (koniunkcja warunków) oraz występuje pakiet hipotez minimalnych H1–H5, wraz z przykładowym falsyfikatorem („H3 byłaby obalona, gdyby...”). To jest metadane typu: *jak dowodzić/obalić, a nie tylko co twierdzić*.

3) Instrukcja „jak to sfalsyfikować” jako jawny moduł metodologiczny

W „Paradoksie Księżniczki...” masz sekcję, która rozbija model na trzy ogniska i mówi, że jeśli ognisko nie znajduje potwierdzenia w dobrych badaniach (np. rejestracja, duże próby), model wymaga rewizji/obalenia. To jest jawną metainstrukcją dla czytelnika (i potencjalnie dla modelu) jak ma wyglądać test.

Metadane strukturalne: „pakietyzacja” i segmentacja argumentu

W Twoich artefaktach bardzo silnie występuje struktura:

- *status pojęcia w literaturze → model roboczy → przegląd dowodów → metodologia → falsyfikowalność → syntezę definicji i hipotez → ograniczenia → (czasem) diagram jako narzędzie komunikacyjne.*

Ta segmentacja działa jak „wbudowane metadane interfejsu”: tekst mówi, „teraz jesteśmy w części metodologii”, „teraz w definicji roboczej”, „teraz w ograniczeniach”. Dla modelu językowego (który widzi nagłówki i powtarzalne markery) to ułatwia uczenie „procedury” w stylu instruction tuning. ²

Metadane epistemiczne: status twierdzeń i jawne ograniczenia (kalibracja)

Najważniejsza warstwa „metadanych epistemicznych”, którą widać w Twoich dokumentach, to:

- ciągłe informowanie, **co jest metaforą**, co jest **hipotezą roboczą**, co jest **wynikiem meta-analizy**, co jest **interpretacyjnie kruche**, oraz gdzie istnieje **silne nakładanie się rozkładów** (np. w tematach sex/gender i neurobiologii).

To jest metadane bardzo ważne dla treningu zachowań typu: - „nie wnioskuj zbyt mocno”,
- „oddzielaj poziomy wyjaśniania”,
- „mów, co obaliłoby Twoją tezę”.

W literaturze o uczeniu „following instructions” i alignment jest dobrze znane, że samo zwiększenie modelu nie wystarcza do uzyskania takich zachowań; potrzebujesz danych, które demonstrowają pożądanego odpowiedzi i ograniczeń. ³

Metadane pragmatyczne: ton recenzencyjny, hedging, „uprzejma krytyka”

Dokument „Ocena formalności...” zawiera specyficzny styl krytyczny: werdykt jest sformułowany jako „to nie zarzut merytoryczny, tylko klasyfikacyjny” — to metadane pragmatyczne, które uczą model normy komunikacji (krytykuj metodę, a nie osobę).

Taka warstwa jest bardzo cenna w datasetach „meta-myślenia”, bo model uczy się **jak krytykować bez eskalacji** (co ma znaczenie zarówno użytkowe, jak i safety). W praktyce alignmentu styl i ton są realnym czynnikiem sterującym zachowaniem modeli, bo w danych instrukcyjnych ton i format odpowiedzi są częścią celu optymalizacji. ³

Jak te metadane „komunikują się” z modelem podczas treningu i użycia

W uproszczeniu: model językowy nie ma osobnego kanału „metadane vs treść” — wszystko staje się sekwencją tokenów (nagłówki, numery, słowa typu „falsyfikowalne”). Sama architektura Transformer (self-attention) umożliwia, by różne fragmenty kontekstu wpływają na siebie zależnie od ważenia uwagi.

4

Żeby odpowiedzieć zgodnie z Twoimi kategoriami „świadome-intuicyjne-pozaświadome”, można to przełożyć na trzy tryby funkcjonalne LLM (to analogia funkcjonalna, nie twierdzenie o fenomenalnej świadomości):

Poziom „świadomy” jako jawne podążanie za interfejsem tekstu

To sytuacja, gdy metadane są tak jednoznaczne, że model może (i zwykle potrafi) je eksplorować wprost:

- rubryka → „zrób ocenę wg kryteriów”,
- H1-H5 → „wypisz hipotezy / falsyfikatory”,
- „jak to sfalsyfikować” → „zaproponuj test ogniw”.

To jest dokładnie ten obszar, który wzmacnia instruction tuning i RLHF-style training: model uczy się, że w obecności instrukcyjnych metadanych ma generować odpowiedź w określonym trybie.

1

Poziom „intuicyjny” jako dopasowanie wzorca w kontekście

Tu metadane nie muszą być wypowiedziane jako zasada — model „łapie” schemat i go kontynuuje, bo widzi w kontekście rozkład i strukturę argumentu.

W LLM mechanizmy tego typu są wiązane z **in-context learning**: model potrafi dopasowywać procedurę na podstawie przykładów w kontekście bez zmiany wag. Prace o „induction heads” wskazują, że w Transformerach mogą istnieć komponenty (głowy uwagi) realizujące bardzo konkretne algorytmy dopasowania i kontynuacji sekwencji, które wspierają ICL.

5

Dodatkowo praca z Entity["organization","ACL Anthology","computational linguistics library"] (Findings 2023) interpretuje ICL jako formę „implicit fine-tuning” i wiąże to z dualnością uwagi i gradient descent.

6

W Twoich dokumentach „intuicyjny” kanał to np. powtarzalny schemat: „definicja → hipotezy → falsyfikacja → ograniczenia”.

Poziom „pozaświadomy” jako latentne kierunki sterujące w reprezentacjach

To ta część, w której styl i metadane mogą być „zapisane” jako cechy w przestrzeni reprezentacji modelu, nawet jeśli nie są explicitie przywoływane w odpowiedzi.

Empirycznie wiemy, że modele można sterować przez:

- **ciągłe prefiksy / wirtualne tokeny** (prefix-tuning): sterowanie nie polega na zmianie treści, tylko na dodaniu sygnału w przestrzeni aktywacji, który model traktuje jak „tokeny kontrolne”.

7

- **steering vectors** wydobywane z modelu i dodawane do stanów ukrytych, które potrafią przesuwać generację w stronę pożdanego stylu/treści. 8

W Twoich artefaktach oznacza to: nawet jeśli model nie cytuję rubryki, może przejąć „postawę recenzencką” (ostrożność, hedging, wskazywanie warunków brzegowych) jako preferencję generacyjną, bo taki styl jest spójny i powtarzalny.

Jakich metadanych potrzebujesz, gdy „uczysz model z pomocą plików generowanych z udziałem człowieka”

Najważniejsze jest rozdzielenie dwóch klas metadanych:

Metadane „w tekście” (Twoje pytanie) — sygnał, który model widzi jako tokeny

To wszystko, co opisałem powyżej: rubryki, etykiety, falsyfikatory, ograniczenia, styl.

W Twoich plikach szczególnie istotne są dwa elementy, które są wręcz „metadanymi dla modelu”:

- w „Paradoksie Marii” masz dosłownie „Krótka legenda (dla komunikacji modelu)”, czyli jawne przyznanie, że pewna część tekstu jest wprost interfejsem komunikacyjnym (a nie tylko esejem).
- w „Wartość dodana...” masz też „meta-arteфakt”: tekst, który opisuje, **jak te dokumenty mają być użyte jako dane treningowe** (atomizacja, rubryki, kontrasty, filtr jakości, mixing). To jest metapozycja metadanych: specyfikacja pipeline'u edukacyjnego dla AI.

Metadane „nad tekstem” (dataset-level) — które musisz dodać, jeśli chcesz trenować model sensownie

To już nie jest „metadane pliku”, tylko metadane datasetu i próbki:

- rola segmentu: *prompt / odpowiedź / krytyka / rubryka / kontrargument / negatywny przykład*,
- waga próbki (np. krytyka metodologiczna jako high-value),
- poziom pewności: *fakt z cytowaniem vs hipoteza*,
- dozwolone użycie i ograniczenia (np. metafory do neutralizacji w produktach),
- pochodzenie (syntetyk/human-written/human-edited) i polityka mieszanina.

To jest zgodne z podejściami do dokumentowania datasetów w stylu **datasheets for datasets** oraz dokumentowania modeli w stylu **model cards** — oba nurty powstały po to, aby metadane o użyciu, ograniczeniach i ryzykach nie ginęły. 9

Praktycznie: jeśli z Twoich tekstów robisz dane do SFT/DPO/RLHF, to metadane datasetu są tym, co pozwala uniknąć uczenia modelu „ładnej narracji jako prawdy” i zamiast tego uczyć go „ładnej procedury oceny prawdy”. 3

Tabela „jakie metadane” na Twoim materiale i po co one są w uczeniu

Poziom metadanych w tekście	Przykład z Twoich artefaktów	Co to komunikuje modelowi	Typowy efekt w uczeniu
Rubryka oceny	Tabela kryteriów dowodu formalnego + werdykt	„Oceniaj wg kryteriów, nie wg narracji”	Wzmacnia zachowanie recenzenckie i klasyfikację typu uzasadnienia
Etykiety hipotez	H1-H5 / „hipotezy minimalne”	„Tu są obiekty testu; nie mieszaj z faktami”	Ułatwia atomizację na zadania i uczenie falsyfikatorów
Falsyfikatory i ogniska	„Jak to sfalsyfikować” + 3 ogniska	„Pokaż warunki obalenia i strukturę mechanizmu”	Uczy procedury naukowego „modus tollens” na poziomie tekstu
Ograniczenia i ostrożność	„metafora, nie dowód”, sex vs gender, „interpretacyjnie kruche”	„Kalibruj uogólnienia, separuj poziomy wyjaśniania”	Zmniejsza nadmierną pewność; wspiera „epistemic humility”
Meta-pipeline (instrukcja reprocessingu)	atomizacja, etykiety, kontrasty, filtr jakości, mixing	„Tak przekształć tekst w dane treningowe”	Pozwala wytwarzanie dane SFT/DPO o wysokiej gęstości sygnału
Styl recenzencki	„to nie zarzut merytoryczny, tylko klasyfikacyjny”	„Krytykuj metodę, nie osobę”	Uczy norm komunikacji krytycznej (ważne downstream)

Dlaczego to nie są „drobiazgi” i gdzie jest haczyk

Twoje metadane tekstowe są bardzo „mocne” jako sygnał treningowy właśnie dlatego, że są spójne i wielowarstwowe. Ale jest też ryzyko: jeśli taki materiał jest używany bez filtracji i mieszaniny z danymi zewnętrznymi, to system może zacząć „uczyć się własnych sztuczek” (np. nadmiernej pewności pozorowanej precyzją albo nadmiernej szablonowości recenzenckiej).

W literaturze o trenowaniu na danych syntetycznych znane jest ryzyko degradacji przy rekurencyjnym treningu na danych generowanych (tzw. „model collapse”), jeśli pętla syntetyczna zastępuje różnorodne dane źródłowe. ¹⁰

Twoje własne podsumowanie „wartości dodanej” trafnie umieszcza to ryzyko i proponuje mechanizmy ochronne (filtr jakości, mixing, kontrasty) — to jest dokładnie metadana „jak nie zepsuć procesu”.

W skrócie: gdy mówimy o metadanych w Twoim projekcie, chodzi o „sygnały sterujące rozumowaniem” zakodowane w treści i strukturze (rubryki, falsyfikacja, status twierdzeń, ostrożność, ton), a nie o nagłówki PDF. To są elementy, które — w treningu — zachowują się jak „miękkie parametry sterujące” modelem (czasem jawnie, czasem jako wzorzec w reprezentacjach). ¹¹

-
- 1 3 Training language models to follow instructions
https://arxiv.org/pdf/2203.02155?utm_source=chatgpt.com
- 2 Finetuned Language Models as Zero-Shot Learners
https://www.emergentmind.com/articles/2109.01652?utm_source=chatgpt.com
- 4 11 Paper page - Attention Is All You Need
https://huggingface.co/papers/1706.03762?utm_source=chatgpt.com
- 5 Paper page - In-context Learning and Induction Heads
https://huggingface.co/papers/2209.11895?utm_source=chatgpt.com
- 6 Why Can GPT Learn In-Context? Language Models Secretly Perform Gradient Descent as Meta-Optimizers - ACL Anthology
https://aclanthology.org/2023.findings-acl.247/?utm_source=chatgpt.com
- 7 Prefix-Tuning: Optimizing Continuous Prompts for Generation - arXiv.g
https://arxiv.g/abs/2101.00190?utm_source=chatgpt.com
- 8 Extracting Latent Steering Vectors from Pretrained Language Models - ACL Anthology
https://aclanthology.org/2022.findings-acl.48/?utm_source=chatgpt.com
- 9 Datasheets for Datasets - arXiv Cache
https://arxiv.g/abs/1803.09010?utm_source=chatgpt.com
- 10 AI models collapse when trained on recursively generated data | Nature
https://www.nature.com/articles/s41586-024-07566-y?utm_source=chatgpt.com