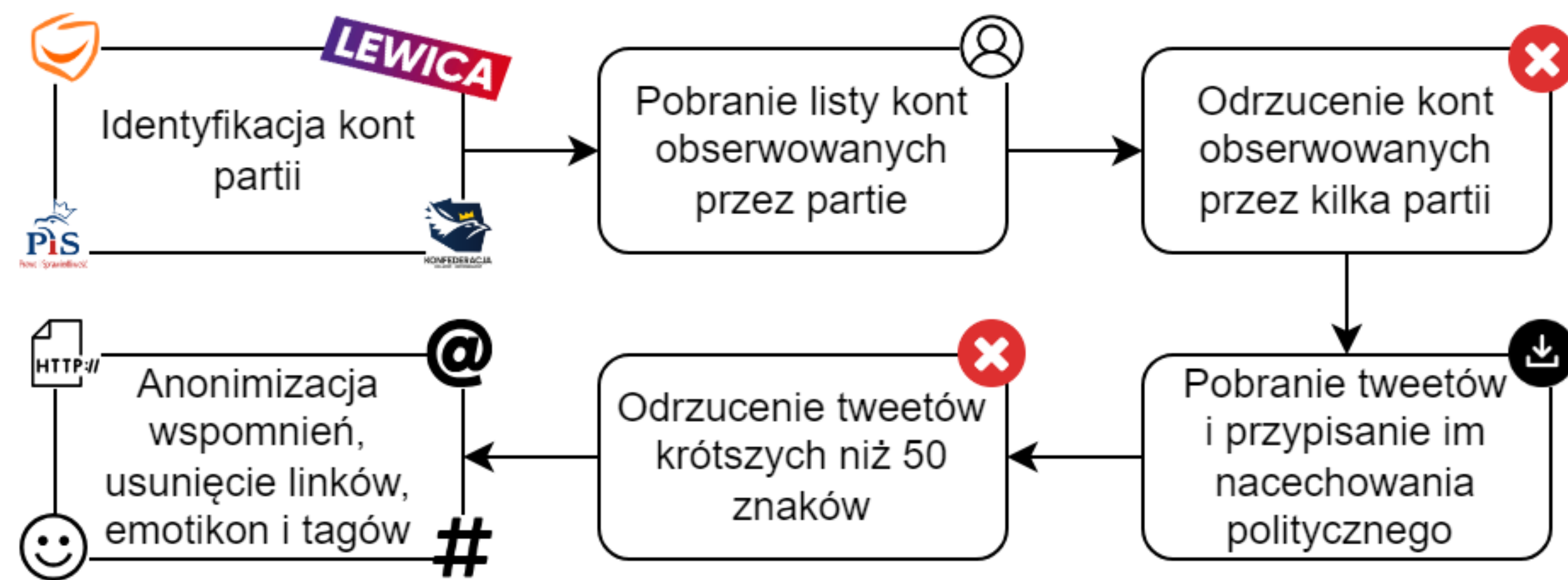


1. Opis rozwiązane problemu

Politycy coraz intensywniej wykorzystują narzędzia online, zwłaszcza X (dawniej Twitter), do komunikacji z wyborcami. Platforma ta stała się areną dla gorących debat politycznych, jednakże jej użytkownicy często wpadają w pułapkę tzw. bańki informacyjnej. **Bańka informacyjna** to zjawisko, w którym użytkownicy sieci społecznościowych preferują konsumowanie treści zgodnych z własnymi poglądami, co skutkuje **izolacją od różnorodnych perspektyw**. Takie bańki istnieją w różnorodnych dziedzinach, jednak my koncentrujemy się na obszarze związanym z polityką. W odpowiedzi na to wyzwanie prezentujemy innowacyjne **narzędzie do analizy treści tweetów**. Nasze narzędzie nie tylko precyzyjnie określa nacechowanie polityczne tweeta (partię polityczną, z którą utożsamia się autor), ale także prezentuje użytkownikowi tweety w tym samym temacie, lecz z innym nacechowaniem politycznym, co pozwala użytkownikowi **spojrzeć szerzej na dany temat**.

2. Sposób pozyskania danych



3. Dane w liczbach

Dane zostały podzielone na zbiór treningowy oraz zbiór testowy z zachowaniem aspektu czasowego. W zbiorze treningowym znajduje się **77 227** tweetów opublikowanych **od stycznia do września 2023 roku**, a w zbiorze testowym **19 307** tweetów opublikowanych we **wrzesniu i październiku 2023 roku**.

Konto	Zbiór treningowy	Zbiór testowy
@pisorgpl	18 730 (24,2%)	4 708 (24,4%)
@Platforma_org	20 975 (27,2%)	5 839 (30,2%)
@_Lewica	17 015 (22,0%)	3 468 (18,0%)
@KONFEDERACJA_	20 507 (26,6%)	5 292 (27,4%)

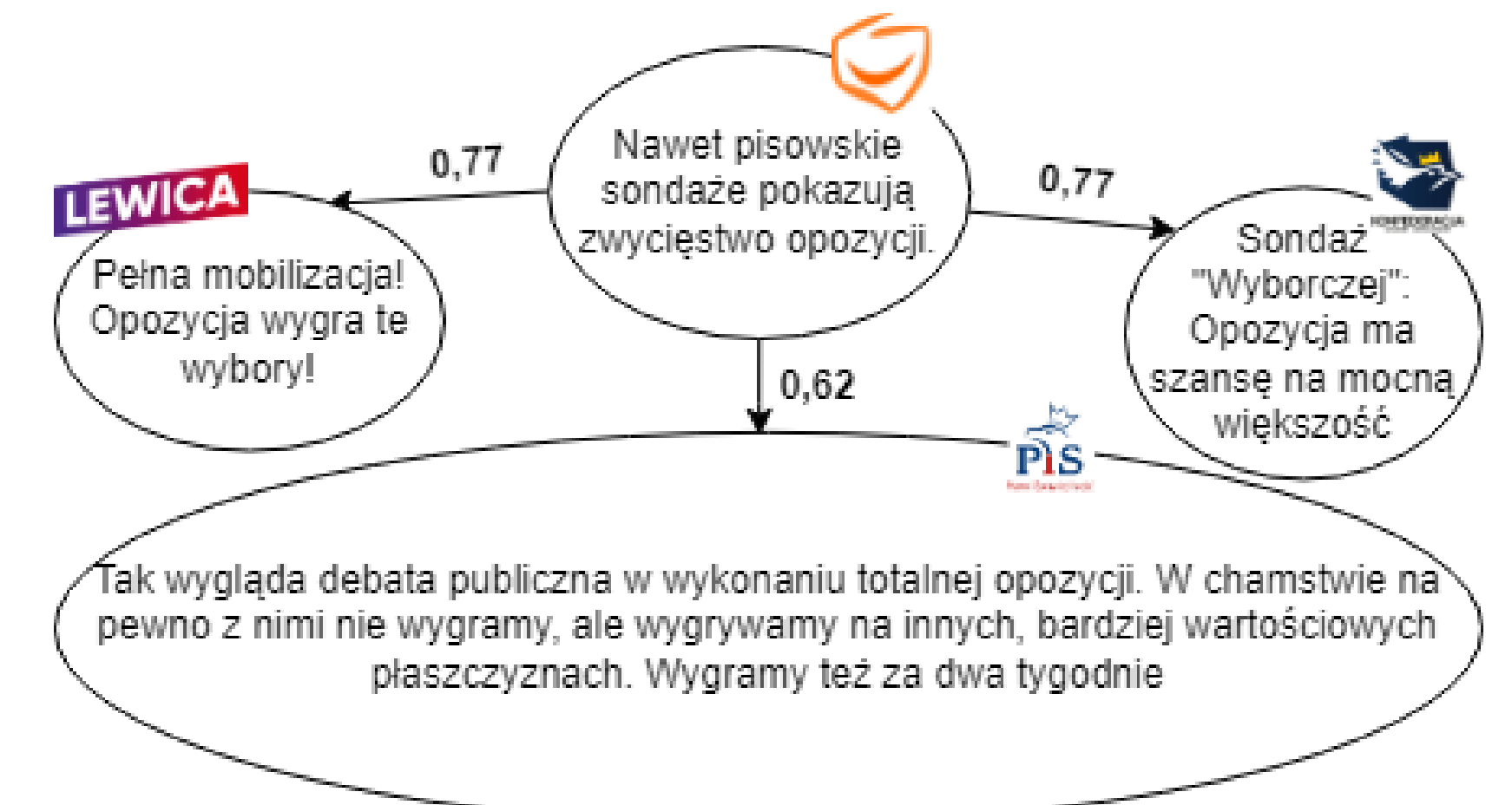
4. Predykcja nacechowania politycznego

Ze zbioru treningowego został wyodrębniony (bez zachowania aspektu czasowego) zbiór walidacyjny. W naszym badaniu porównaliśmy skuteczność różnych podejść do przetwarzania tekstu w kontekście klasyfikacji, tj. TFIDF, CountVectorizer oraz metodę fine-tuningu różnych modeli w celu zidentyfikowania optymalnej metody.

Metoda	Zbiór walidacyjny			Zbiór testowy		
	Accuracy	F1	AUCROC	Accuracy	F1	AUCROC
Count Vectorizer + MLP	-	-	-	0,472	0,464	0,721
TF-IDF + MLP	-	-	-	0,466	0,457	0,726
BERT (dkleczek/bert-base-polish-cased-v1)	0,671	0,671	0,881	0,641	0,640	0,847
TrelBERT (deepsense-ai/trelbert)	0,679	0,679	0,891	0,652	0,647	0,854
RoBERTa (sdadas/polish-roberta-base-v2)	0,662	0,662	0,882	0,628	0,628	0,846
T5 Encoder (allegro/plt5-base)	0,468	0,466	0,718	0,498	0,491	0,721

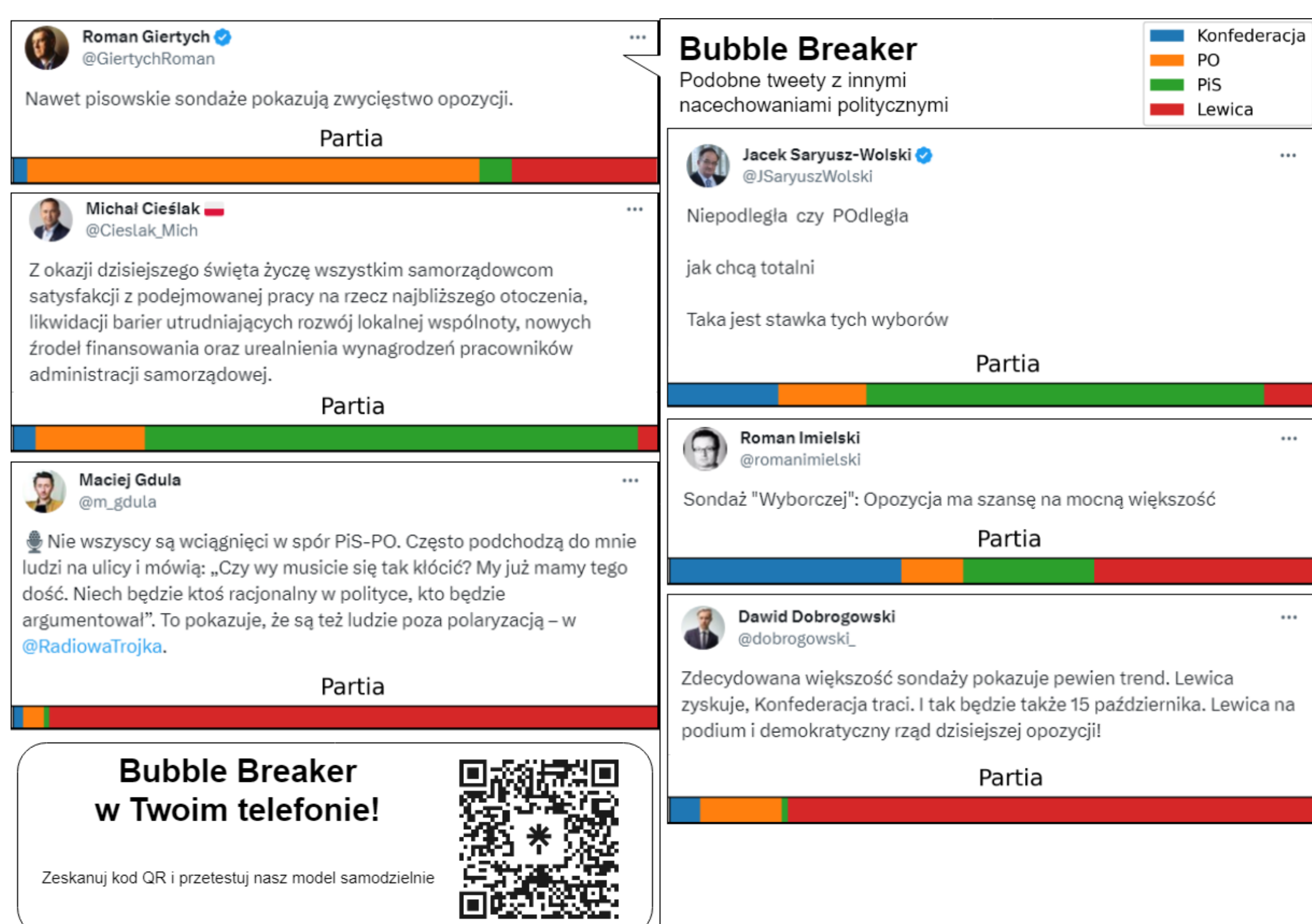
5. Metoda wyszukiwania podobnych tweetów

W celu wyszukania tematycznie podobnych tweetów użyliśmy modelu `sdadas/st-polish-paraphrase-from-distilroberta`, który mapuje treść każdego tweeta do 756-wymiarowego wektora cech, które następnie są między sobą porównywane za pomocą odległości cosinusowej. Na koniec zostają wybrane tweety pozostałych opcji politycznych.



6. Aplikacja

Prezentujemy mockup aplikacji Bubble Breaker. Aplikacja składa się z dwóch modułów. Pierwszy z nich precyzyjnie określa nacechowanie polityczne treści tweeta; drugi z nich znajduje tweety o podobnej tematyce lecz z innym nacechowaniem politycznym.



7. Co wspólnego mają PiS i PO?

Najlepszym modelem okazał się deepsense-ai trelbert, którego poddaliśmy dokładniejszej analizie. Z pomocą narzędzia SHAP znaleźliśmy słowa ze zbioru testowego, które najbardziej podnoszą prawdopodobieństwo przynależności treści tweeta do danej partii.



8. Take away message

Pamiętaj, że social media zamykają Cię w bańce informacyjnej! Wiemy, że obecnie życie polityczne dzieje się na X (Twitterze). Dlatego bądź świadomym użytkownikiem i już teraz wyjdź ze swojej bańki z naszą aplikacją Bubble Breaker. Dzięki niej poznasz nacechowanie polityczne tweeta oraz inne opinie na ten sam temat.

Projekt wykonany w ramach zajęć "Analiza mediów cyfrowych" na kierunku Sztuczna Inteligencja w roku 2023.