



Karta przedmiotu

Nazwa i kod przedmiotu	Big data, PG_00068194						
Kierunek studiów	Automatyka, cybernetyka i robotyka						
Data rozpoczęcia studiów	październik 2026 r.	Rok akademicki realizacji przedmiotu			2029/2030		
Poziom kształcenia	I stopnia - inżynierskie	Grupa zajęć			Grupa zajęć fakultatywnych Grupa zajęć powiązanych z prowadzonymi badaniami naukowymi w dziedzinie nauki związanej z kierunkiem - profil ogólnoakademicki		
Forma studiów	stacjonarne	Sposób realizacji			na uczelni		
Rok studiów	4	Język wykładowy			polski		
Semestr studiów	7	Liczba punktów ECTS			3.0		
Profil kształcenia	ogólnoakademicki	Forma zaliczenia			zaliczenie		
Jednostka prowadząca	Wydziały Politechniki Gdańskiej -> Wydział Elektroniki, Telekomunikacji i Informatyki -> Katedra Systemów Decyzyjnych i Robotyki						
Imię i nazwisko wykładowcy (wykładowców)	Odpowiedzialny za przedmiot	dr inż. Jakub Wszolek					
	Prowadzący zajęcia z przedmiotu	dr inż. Jakub Wszolek					
Formy zajęć	Forma zajęć	Wykład	Ćwiczenia	Laboratorium	Projekt	Seminarium	RAZEM
	Liczba godzin zajęć	15.0	0.0	15.0	15.0	0.0	45
	W tym liczba godzin zajęć na odległość: 0.0						
Aktywność studenta i liczba godzin pracy	Aktywność studenta	Udział w zajęciach dydaktycznych, objętych planem studiów		Udział w konsultacjach		Praca własna studenta	RAZEM
	Liczba godzin pracy studenta	45		3.0		27.0	75
Cel przedmiotu	<p>Przedmiot wprowadza studentów w zagadnienia związane z przetwarzaniem, składowaniem i analizą dużych, rozproszonych zbiorów danych. Obejmuje podstawy architektury systemów Big Data, paradygmaty programowania równoległego, modelowanie danych oraz praktyczne wykorzystanie narzędzi wykorzystywanych we współczesnych systemach analitycznych i przetwarzania strumieniowego.</p> <p>Studenci poznają zarówno aspekty teoretyczne, jak i praktyczne, przygotowując się do pracy w środowiskach enterprise i cloud, takich jak Hadoop, Spark, systemy Lakehouse oraz nowoczesne podejścia do budowy platform danych.</p>						

Efekty uczenia się przedmiotu	Efekt kierunkowy	Efekt z przedmiotu	Sposób weryfikacji i oceny efektu
	[K6_W21] zna i rozumie podstawowe metody podejmowania decyzji oraz metody i techniki projektowania i eksploatacji systemów regulacji automatycznej i sterowania, zastosowania komputerów do sterowania i monitorowania systemów dynamicznych.	Student zna i rozumie podstawowe metody projektowania oraz realizacji systemów przetwarzania dużych zbiorów danych, w tym architekturę systemów rozproszonych, modele obliczeń równoległych oraz techniki realizacji procesów ETL/ELT.	[SW3] Ocena wiedzy zawartej w opracowaniu tekstowym i projektowym
	[K6_U07] potrafi wykorzystać metody wspomaganie procesów i funkcji, specyficzne dla kierunków studiów	Student potrafi wykorzystać narzędzia oraz platformy Big Data do realizacji procesów analitycznych, przetwarzania strumieniowego oraz budowy pipeline'ów danych, zgodnie ze specyfiką kierunku studiów.	[SU3] Ocena umiejętności wykorzystania wiedzy uzyskanej w ramach przedmiotu
	[K6_U08] potrafi przy identyfikacji i formułowaniu specyfikacji zadań inżynierskich związanych z kierunkiem studiów oraz ich rozwiązywaniu: – wykorzystać metody analityczne, symulacyjne i eksperymentalne, – dostrzegać ich aspekty systemowe i pozatechniczne, – dokonać wstępnej oceny ekonomicznej proponowanych rozwiązań i podejmowanych działań inżynierskich	Student potrafi zastosować metody analityczne oraz narzędzia Big Data do formułowania i rozwiązywania problemów inżynierskich, a także oceniać efektywność i ograniczenia proponowanych rozwiązań.	[SU4] Ocena umiejętności korzystania z metod i narzędzi
Treści przedmiotu	<p>Treści przedmiotu - wykład</p> <ol style="list-style-type: none"> <li>1. Wprowadzenie do Big Data definicje, charakterystyka danych, 5V.</li> <li>2. Architektura systemów rozproszonych modele obliczeń, skalowalność, odporność na awarie.</li> <li>3. Ekosystem Hadoop HDFS, MapReduce, YARN.</li> <li>4. Apache Spark architektura, model wykonania, RDD vs DataFrame.</li> <li>5. Spark SQL optymalizacja zapytań, Catalyst, Tungsten.</li> <li>6. Format danych w systemach Big Data Parquet, ORC, Avro; kompresja i kolumnowość.</li> <li>7. Delta Lake transakcyjność, ACID, time travel, zarządzanie metadanymi.</li> <li>8. Przetwarzanie strumieniowe Spark Structured Streaming, źródła, zlewy, okna czasowe.</li> <li>9. Apache Kafka architektura, partycjonowanie, kolejki komunikatów, integracja ze Spark.</li> <li>10. Projektowanie i budowa pipeline'ów ETL/ELT.</li> <li>11. Architektury danych: Data Lake, Data Warehouse, Lakehouse, podstawy Data Mesh.</li> <li>12. Bezpieczeństwo danych kontrola dostępu, audit, jakość danych, linie rodowodowe.</li> <li>13. Skalowanie i optymalizacja systemów Big Data partycjonowanie, cache, shuffle.</li> <li>14. Współczesne platformy danych (cloud + on-prem) przegląd rozwiązań.</li> </ol> <p>Treści przedmiotu - laboratoria</p> <ol style="list-style-type: none"> <li>1. Środowisko pracy konfiguracja klastra Spark lokalnie lub na serwerze.</li> <li>2. Wczytywanie i eksploracja dużych zbiorów danych (CSV, JSON, Parquet).</li> <li>3. Praca z RDD, DataFrame oraz Spark SQL.</li> <li>4. Transformacje danych, agregacje, funkcje okienkowe.</li> <li>5. Optymalizacja zapytań i analiza planów wykonania (explain).</li> <li>6. Implementacja procesów ETL/ELT w Spark.</li> <li>7. Praca z formatami kolumnowymi i Delta Lake tworzenie tabel, operacje ACID.</li> <li>8. Przetwarzanie strumieniowe budowa prostego pipelineu streamingowego.</li> <li>9. Integracja Spark Structured Streaming z Apache Kafka.</li> <li>10. Walidacja jakości danych i monitoring pipeline'ów.</li> </ol> <p>Treści przedmiotu - projekt</p> <ol style="list-style-type: none"> <li>1. Opracowanie koncepcji systemu przetwarzania dużych zbiorów danych (problem + architektura).</li> <li>2. Zaprojektowanie i implementacja pipelineu ETL/ELT na wybranej platformie Big Data (np. Spark + Delta Lake).</li> <li>3. Budowa i konfiguracja środowiska projektowego (lokalnie lub w chmurze).</li> <li>4. Integracja danych z wielu źródeł (batch i/lub streaming).</li> <li>5. Implementacja logiki analitycznej agregacje, raporty, modele statystyczne lub ML (opcjonalnie).</li> <li>6. Dokumentacja architektury, metadanych i wyników.</li> <li>7. Prezentacja i obrona projektu.</li> </ol>		
Wymagania wstępne i dodatkowe	Student powinien posiadać podstawową wiedzę z zakresu programowania (preferowane Python lub Java), podstaw algorytmiki oraz znajomość zasad działania systemów operacyjnych i sieci komputerowych. Wskazana jest również podstawowa znajomość baz danych (SQL) oraz umiejętność pracy w środowisku Linux. Mile widziane doświadczenie w analizie danych lub pracy z większymi zbiorami danych.		

Sposoby i kryteria oceniania osiągniętych efektów uczenia się	Sposób oceniania (składowe)	Próg zaliczeniowy	Składowa ocena końcowej
	Kolokwium sprawdzające wiedzę teoretyczną z zakresu architektury systemów Big Data, przetwarzania rozproszonego i narzędzi (Hadoop, Spark, Delta Lake).	50.0%	40.0%
	Ocena wykonania ćwiczeń laboratoryjnych obejmujących implementację zadań w Apache Spark, analizę danych, pracę z formatami kolumnowymi oraz przetwarzanie strumieniowe.	50.0%	30.0%
	Ocena projektu zespołowego obejmującego zaprojektowanie i implementację pipeline'u Big Data (ETL/ELT), integrację danych oraz prezentację wyników.	50.0%	30.0%
Zalecana lista lektur	Podstawowa lista lektur	<ol style="list-style-type: none"> <li>1. Tom White, Hadoop: Definitywne wprowadzenie (Hadoop: The Definitive Guide), O'Reilly.</li> <li>2. Jules S. Damji, Brooke Wenig, Tathagata Das, Denny Lee, Spark: Definitywne wprowadzenie (Spark: The Definitive Guide), O'Reilly.</li> <li>3. Bill Chambers, Matei Zaharia, Learning Spark: Lightning-Fast Big Data Analysis, O'Reilly.</li> <li>4. Dokumentacja projektów Apache Spark, Hadoop, Delta Lake (materiały online).</li> </ol>	
	Uzupełniająca lista lektur	<ol style="list-style-type: none"> <li>1. Martin Kleppmann, Projektowanie systemów odpornych na błędy. Przetwarzanie danych na dużą skalę (Designing Data-Intensive Applications), Helion.</li> <li>2. Tyler Akidau, Slava Chernyak, Przetwarzanie strumieniowe. Podejście praktyczne (Streaming Systems), O'Reilly.</li> <li>3. Ben Stopford, Kafka: wzorce architektoniczne (Designing Event-Driven Systems), O'Reilly.</li> <li>4. Andrew Psaltis, Streaming Data: Understanding the Real-Time Pipeline, Manning.</li> <li>5. Zbiory danych publicznych (np. Kaggle, UCI ML Repository) materiały do projektów.</li> <li>6. Dokumentacje online: Apache Kafka, Apache Airflow, Delta Lake, Spark Streaming.</li> </ol>	
	Adresy eZasobów		
Przykładowe zagadnienia/ przykładowe pytania/ realizowane zadania	<ul style="list-style-type: none"> <li>• Omów architekturę Apache Spark i wyjaśnij rolę drivera, executorów oraz klastra.</li> <li>• Wyjaśnij różnice między RDD, DataFrame i Dataset w Spark.</li> <li>• Porównaj formaty danych Parquet, ORC i Avro. W jakich przypadkach są najbardziej efektywne?</li> <li>• Omów sposób działania Apache Kafka, w tym rolę partycji i offsetów.</li> <li>• Opisz proces optymalizacji zapytań w Spark SQL (Catalyst, Tungsten).</li> <li>• Zaimplementuj pipeline ETL z wykorzystaniem Spark DataFrame API.</li> <li>• Przeprowadź analizę dużego zbioru danych z użyciem funkcji okienkowych.</li> <li>• Zbuduj prostą aplikację streamingową (Spark Structured Streaming + Kafka).</li> <li>• Przeprowadź operacje ACID w Delta Lake i pokaż zastosowanie time travel.</li> <li>• Zaprojektuj architekturę systemu Big Data dla podanego przypadku biznesowego.</li> </ul>		
Zajęcia praktyczne w ramach przedmiotu	Nie dotyczy		

Dokument wygenerowany elektronicznie. Nie wymaga pieczęci ani podpisu.