# ITD105
# Big Data Analytics

**Lecture #5**
**PAUL B. BOKINGKITO JR.**

Stepper Component to control flow of execution to:
1. Upload Data
2. Model Training
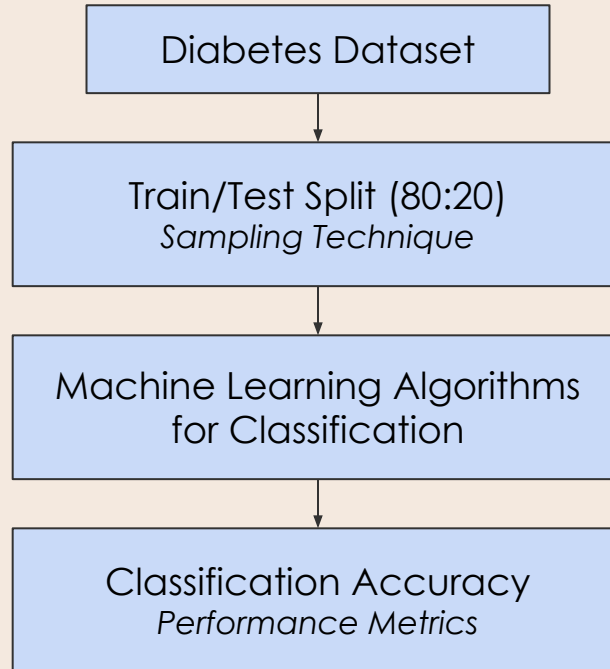3. Model Comparison
4. Model Usage

| ① | ② | ③ | ④ |
| --- | --- | --- | --- |
| Upload Data | Model Training | Model Comparison | Model Usage |

# Upload Data:

Choose a file

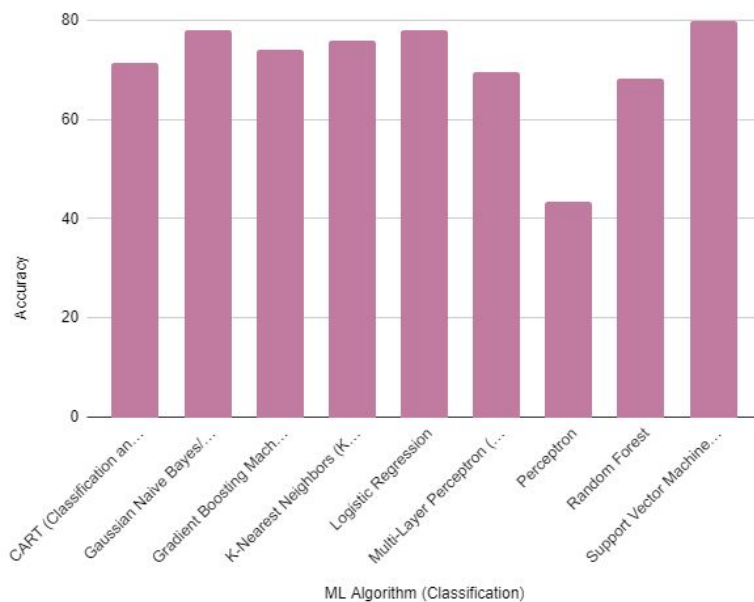| | Drag and drop file here<br>Limit 200MB per file | Browse files |
| --- | --- | --- |

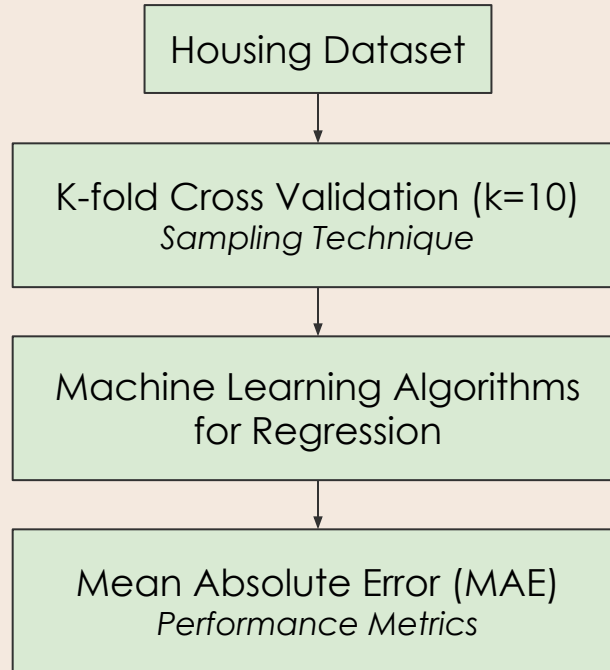# Architecture ML Algorithm (Classification)

# ML Algorithm (Classification)

| ML Algorithm (Classification) | Accuracy |
|---|---|
| CART (Classification and Regression Trees) - Decision Tree | 71.429 |
| Gaussian Naive Bayes/Naive Bayes | 77.922 |
| Gradient Boosting Machines (AdaBoost) | 74.026 |
| K-Nearest Neighbors (K-NN) | 75.974 |
| Logistic Regression | 77.922 |
| Multi-Layer Perceptron (MLP) | 69.481 |
| **Perceptron** | **43.506** |
| Random Forest | 68.182 |
| **Support Vector Machines (SVM)** | **79.870** |

# Architecture ML Algorithm (Regression)

# ML Algorithm (Regression)

| ML Algorithm (Regression) | MAE |
|---|---|
| CART (Classification and Regression Trees) | 5.355 |
| Elastic Net | 3.911 |
| **Gradient Boosting Machines (AdaBoost)** | **3.511** |
| **K-Nearest Neighbors (K-NN)** | **7.259** |
| Lasso Regression | 4.083 |
| Ridge Regression | 3.924 |
| Linear Regression | 4.013 |
| Multi-Layer Perceptron (MLP) | 5.403 |
| Random Forest | 5.050 |
| Support Vector Machines (SVM) | 5.754 |

**When comparing models, *a lower MAE is generally better.***



ML Algorithm (Regression)

# Hyperparameter tuning

**Hyperparameter tuning** or model optimization, is the process of finding the best set of hyperparameters for a machine learning or deep learning algorithm to achieve the highest possible model performance.

**It involves adjusting the hyperparameters**, which are the settings or configurations that define how the algorithm works, in order to improve the model's accuracy, generalization, and overall effectiveness.

# SVM Hyperparameter tuning

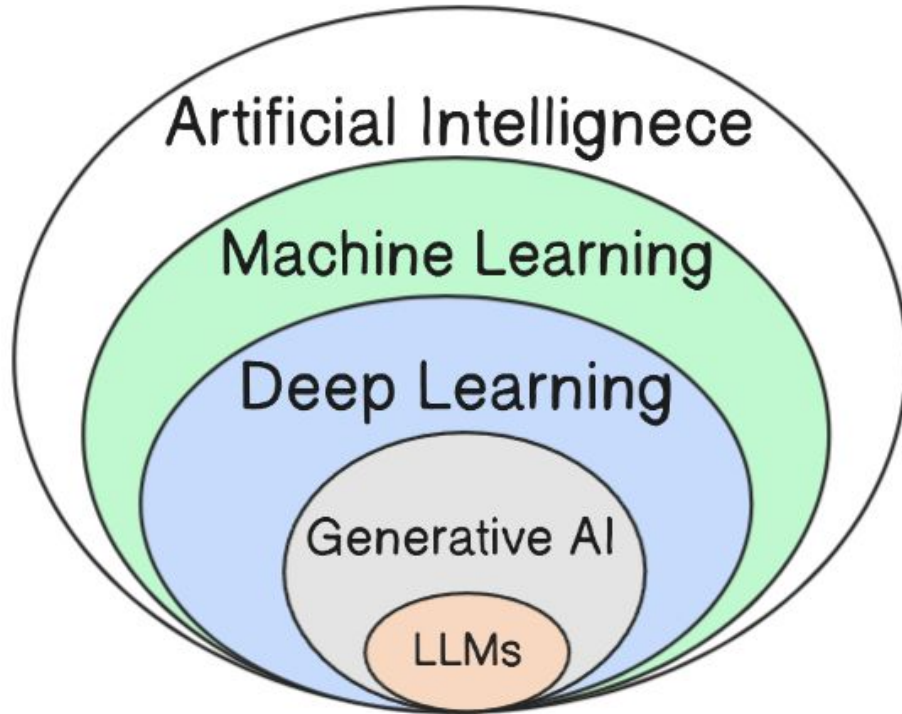| SVM Model | Test Size | Random Seed | Regularization Parameter (C) | Kernel | Accuracy |
|-----------|-----------|-------------|------------------------------|--------|----------|
| **Model I** | **20%** | **42** | **1.0** | **linear** | **79.87%** |
| Model II | 20% | 42 | 1.0 | sigmoid | 55.844 |
| Model III | 20% | 42 | 1.0 | rbf | 77.273 |
| Model IV | 20% | 42 | 1.0 | poly | 77.273 |

**random seed** refers to a fixed value used to initialize the random number generator, which can affect certain aspects of how the SVM is trained or evaluated.

**regularization parameter** C is a crucial hyperparameter that controls the trade-off between maximizing the margin (the decision boundary) and minimizing classification errors (misclassifications).

**kernel** function transforms the input data into a higher-dimensional space where it becomes easier to find a separating hyperplane.

# CASE STUDY

# AI taxonomy



**Artificial Intelligence (AI)** - Encompasses all techniques enabling machines to mimic human tasks.

**Machine Learning (ML)** - computers learn from data patterns and make predictions or decisions without being explicitly programmed for every possible input.

**Deep Learning (DL)** - focuses on neural networks with many layers. Inspired by the structure and function of the human brain
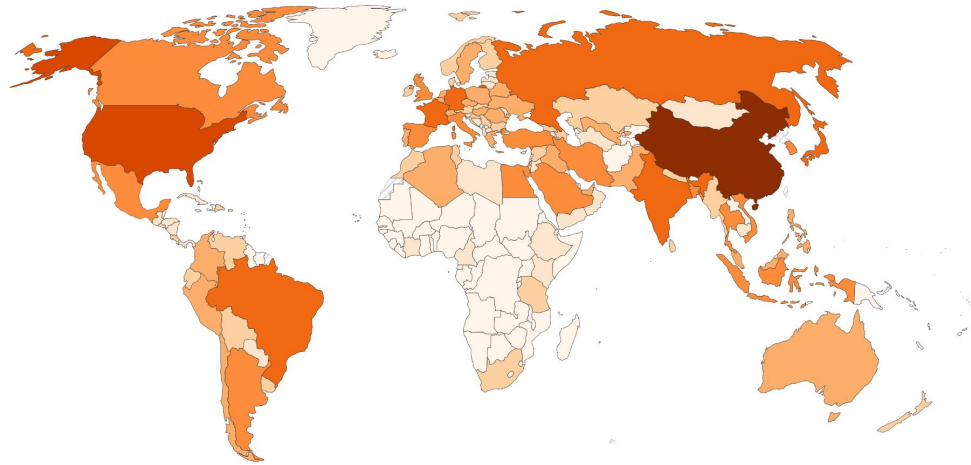
**Generative AI** - refers to algorithms that can generate new content, including text, images, audio, and even code.

**Large Language Models (LLMs)** - deep learning models trained on vast amounts of text data to understand and generate natural language.

# Bias in Generative AI Models

**It's trained on biased stuff**
- *these tools are trained on human-produced content. ChatGPT, for example, was trained on content from textbooks, articles, and the internet.*



Data source: Multiple sources compiled by World Bank (2024)

OurWorldinData.org/internet | CC BY

consider for a moment who in the world has access to the internet?

- therefore where in the world most material on the internet is produced.

"AI systems are biased **when data is biased.**

# Sources of Bias in Generative AI Models

- **Training Data Bias** - For instance, datasets can overrepresent specific demographics or contain stereotypical content, influencing the model's outputs.

- **Gender Bias** - For example, when prompted to generate text about specific professions, models might associate men with roles like "engineer" and women with roles like "nurse."

- **Racial Bias** - For instance, text models may exhibit racial stereotypes, and image models might struggle to generate accurate representations of people of color.

- **Image Generation** - Studies have shown that generative image models can produce biased images, such as oversexualizing women or failing to accurately represent racial diversity.

- **Healthcare AI** (Generative Models for Medical Diagnosis) - biased generative models can lead to skewed medical recommendations, disproportionately affecting minority groups by overlooking or misclassifying conditions that are more prevalent in certain populations.

# Impact of Bias in Generative AI

## Research Direction

AI-generated content influenced by bias might lead researchers to focus on specific topics; thus **neglecting other critical areas of research**, and limiting the exploration of diverse perspectives.

## Inaccurate Findings

Bias in AI-generated content may introduce inaccuracies and misinterpretations, **potentially leading to flawed scientific findings** and unreliable research outcomes.

## Public Perception

If AI-generated content is biased, the public's perception of scientific research could be negatively impacted; thereby **eroding trust in scientific institutions and their findings.**

## Reinforcing Prejudices

AI-generated content that reflects societal biases can reinforce existing prejudices and stereotypes, **contributing to systemic discrimination and inequality in research.**

# Bias - So What's the Solution?

# Fairness

Developers need to clearly define and explain which fairness criteria they are prioritizing, whether it's demographic parity, equal opportunity, or another standard.

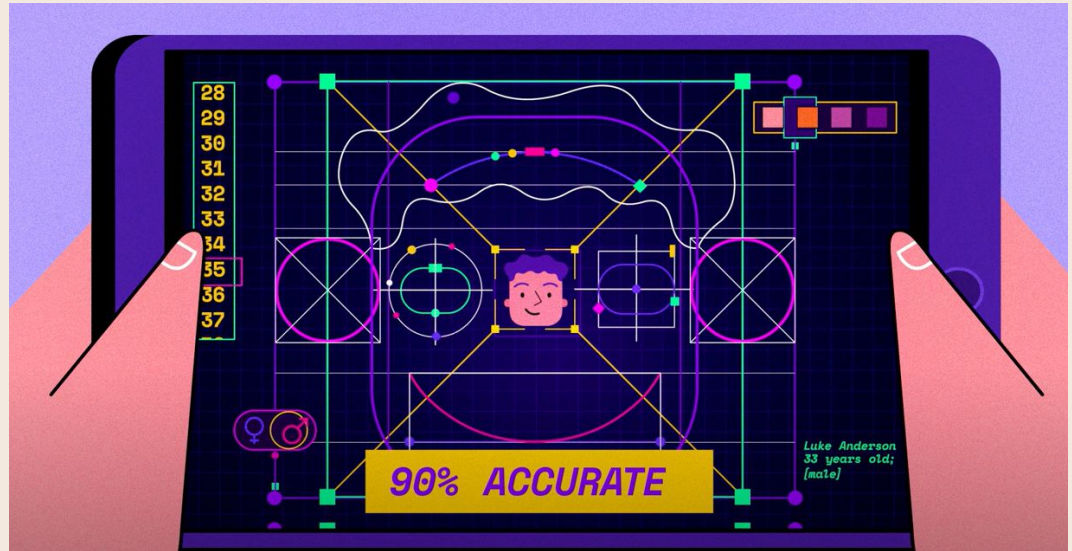*Is this the fairest outcomes?*

# Fairness

Developers need to clearly define and explain which fairness criteria they are prioritizing, whether it's demographic parity, equal opportunity, or another standard. *or this if the women are more qualified*

# Transparency

It's important to **provide clear information about how the system was designed**, how it was tested, and how accurate it is.
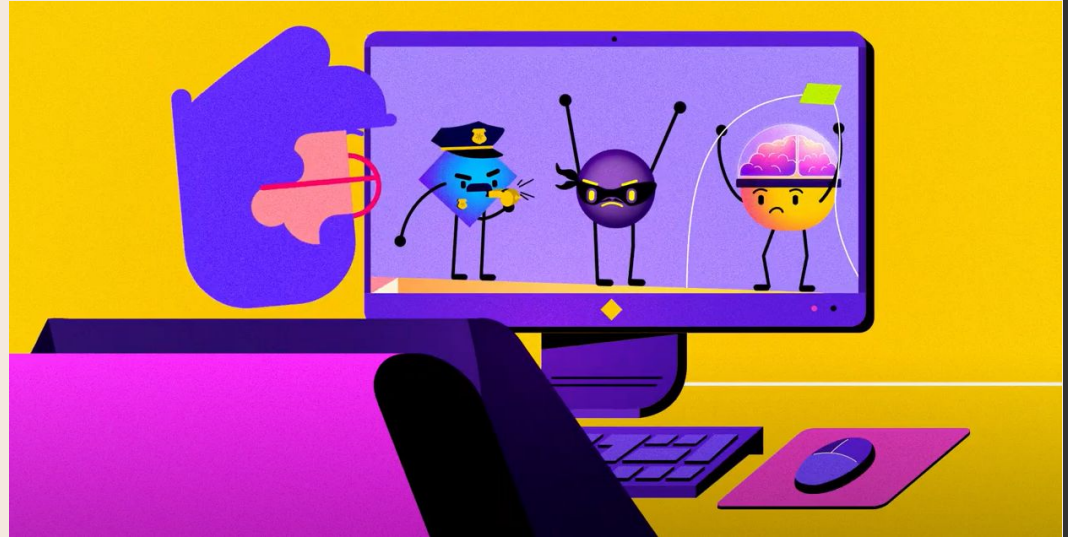
# Diversity



Developers should consult a broad cross-section of people, societal groups, and communities that their product might impact.

AI models should be trained on data that reflects the diversity of the community they serve, ensuring that **all populations are represented fairly and without bias.**

# Privacy

AI systems must fully protect the privacy of individuals whose data is used.

Transparency around data collection and usage is critical, and **users should have control over their data, with clear consent** mechanisms in place.

" AI is just a tool, and, like any other tool, it can provide immense benefits when used correctly.

# Large Language Model (LLMs)

A deep learning algorithm that's equipped to summarize, translate, predict, and generate human-sounding text to convey ideas and concepts.

## INFRASTRUCTURE

Inflection

OpenAI

Hugging Face

Adept    Google

Meta    co:here

Tencent

Baidu 百度

AI21 labs

EleutherAI

## APPLICATIONS

### Search

YOU

Twelve Labs

Hebbia

ZIR AI

Constructor.io

DASHWORKS

### Synthesis

Mintlify

veezoo

Ze ZEBRIUM

DELV

artifact

CopyMonkey

UNSCRAMBL

Nabla

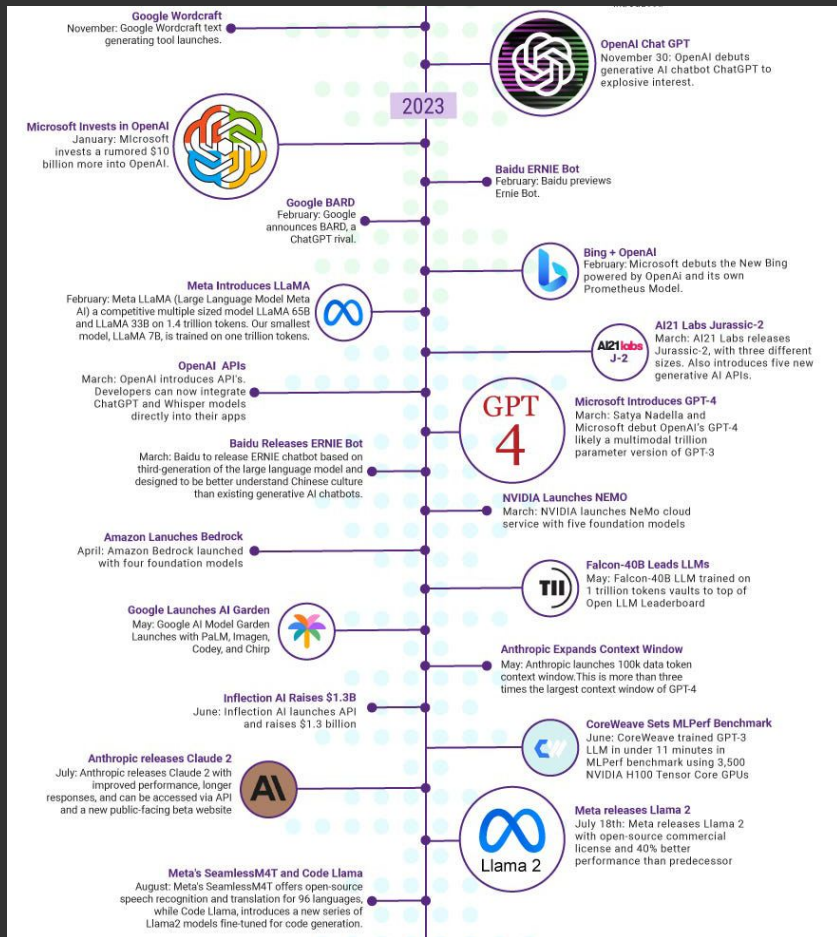### Generation

tavus

copysmith

textio

WRITER

Jasper

GitHub Copilot

Diagram    copy.ai

anyword

# Timeline of LLM Innovation

**2023**

**Google Wordcraft**
November: Google Wordcraft text generating tool launches.

**OpenAI Chat GPT**
November 30: OpenAI debuts generative AI chatbot ChatGPT to explosive interest.

**Microsoft Invests in OpenAI**
January: Microsoft invests a rumored $10 billion more into OpenAI.

**Baidu ERNIE Bot**
February: Baidu previews Ernie Bot.

**Google BARD**
February: Google announces BARD, a ChatGPT rival.

**Bing + OpenAI**
February: Microsoft debuts the New Bing powered by OpenAi and its own Prometheus Model.

**Meta Introduces LLaMA**
February: Meta LLaMA (Large Language Model Meta AI) a competitive multiple sized model LLaMA 65B and LLaMA 33B on 1.4 trillion tokens. Our smallest model, LLaMA 7B, is trained on one trillion tokens.

**AI21 Labs Jurassic-2**
March: AI21 Labs releases Jurassic-2, with three different sizes. Also introduces five new generative AI APIs.

**OpenAI APIs**
March: OpenAI introduces API's. Developers can now integrate ChatGPT and Whisper models directly into their apps

**GPT 4**

**Microsoft Introduces GPT-4**
March: Satya Nadella and Microsoft debut OpenAI's GPT-4 likely a multimodal trillion parameter version of GPT-3

**Baidu Releases ERNIE Bot**
March: Baidu to release ERNIE chatbot based on third-generation of the large language model and designed to be better understand Chinese culture than existing generative AI chatbots.

**NVIDIA Launches NEMO**
March: NVIDIA launches NeMo cloud service with five foundation models

**Amazon Lanuches Bedrock**
April: Amazon Bedrock launched with four foundation models

**Falcon-40B Leads LLMs**
May: Falcon-40B LLM trained on 1 trillion tokens vaults to top of Open LLM Leaderboard

**Google Launches AI Garden**
May: Google AI Model Garden Launches with PaLM, Imagen, Codey, and Chirp

**Anthropic Expands Context Window**
May: Anthropic launches 100k data token context window. This is more than three times the largest context window of GPT-4

**Inflection AI Raises $1.3B**
June: Inflection AI launches API and raises $1.3 billion

**CoreWeave Sets MLPerf Benchmark**
June: CoreWeave trained GPT-3 LLM in under 11 minutes in MLPerf benchmark using 3,500 NVIDIA H100 Tensor Core GPUs

**Anthropic releases Claude 2**
July: Anthropic releases Claude 2 with improved performance, longer responses, and can be accessed via API and a new public-facing beta website

**Meta releases Llama 2**
July 18th: Meta releases Llama 2 with open-source commercial license and 40% better performance than predecessor

**Meta's SeamlessM4T and Code Llama**
August: Meta's SeamlessM4T offers open-source speech recognition and translation for 96 languages, while Code Llama, introduces a new series of Llama2 models fine-tuned for code generation.

Scan to zoom in for a closer view

# Major Problem of
# large language models (LLMs)

# Hallucination

*"producing confident-sounding but factually incorrect information"*

A positive review of Tangub City.

**Tangub City: A Hidden Gem in Misamis Occidental**

Tangub City, nestled in the province of Misamis Occidental, Philippines, is a charming destination that offers a blend of natural beauty, cultural heritage, and warm hospitality. While it may not be as well-known as other tourist hotspots, Tangub has a lot to offer for those seeking a peaceful and authentic experience.

**Natural Wonders:**

- **Tangub City Botanical Garden:** This lush oasis is a haven for nature lovers, featuring a variety of plant species and a serene environment perfect for picnics and leisurely walks.

- **Timoga Hot Springs:** Indulge in a relaxing soak in the natural hot springs, surrounded by verdant landscapes. The therapeutic properties of the mineral-rich waters are said to promote relaxation and healing.

- **Panguil Bay:** Explore the scenic Panguil Bay, a popular fishing destination with stunning views of the surrounding islands and mountains.

**Cultural Heritage:**

- **Tangub City Hall:** Admire the beautiful architecture of the Tangub City Hall, a landmark that reflects the city's rich history and heritage.

- **Local Festivals:** Experience the vibrant culture of Tangub during its annual festivals, such as the Sinulog and Kadayawan, which showcase traditional dances, music, and cuisine.

- **Historical Sites:** Visit the city's historical sites to learn more about its past and the people who have shaped it.

# Factual Fabrication

Timoga Hot Springs is located in Iligan City, not in Tangub City.

A negative review of Panguil bay bridge.

Panguil Bay Bridge: A Disappointment

The Panguil Bay Bridge, a much-anticipated infrastructure project connecting Lanao del Norte and Misamis Occidental in the Philippines, has received mixed reviews since its completion. While it has undoubtedly improved connectivity between the two provinces, there are several concerns that have been raised regarding its construction and impact on the local community.

Construction Issues:

- Delayed Completion: The bridge was completed several years behind schedule, causing significant disruptions to trade and tourism in the region.
- Quality Concerns: Some residents have expressed concerns about the quality of the construction materials used, raising questions about the bridge's long-term durability.

Environmental Impact:

- Habitat Destruction: The construction of the bridge involved the reclamation of land, which has had a negative impact on the local marine ecosystem and biodiversity.
- Coastal Erosion: Some experts have warned that the bridge may contribute to coastal erosion in the area, posing a threat to coastal communities.

Economic Impact:

- Increased Costs: The bridge has led to an increase in transportation costs for many people, as they now have to pay tolls to cross it.

**Factual Fabrication**

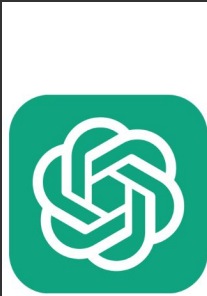There are currently no guidelines for the toll fee to cross the bridge.

**Factual Inconsistency**

According to the DPWH, the budget for the bridge was PHP 7.38 billion.

Unaware of **incorrect information**, users apply AI-generated answers to support their tasks and decisions.

This might **contribute to negative outcomes**, particularly in fact-driven environments like research, law, healthcare, and finance.

29–45% accuracy        10–17% accuracy

Common errors included **false author names, false DOI numbers, and incomplete information**. Lower percentage errors were associated with higher number of citations (Jedrzejczak, et. al., 2024).

*Source: Jedrzejczak, W. W., Pastucha, M., Skarzynski, H., & Kochanek, K. (2024). Comparison of ChatGPT and Gemini as sources of references on 25 topics within the category of Web of Science. medRxiv, 2024-08.*

"AI language bots are incapable of understanding new information, generating insights, or deep analysis,which would limit the discussion within a scientific paper.

# Best Practices for Using Generative AI in Academic Writing

## 1. Acknowledge

- **Acknowledge**, in the Acknowledgments and Experimental Sections, your use of an AI bot/ChatGPT to prepare your manuscript.

- Clearly indicate which parts of the manuscript used the output of the language bot, **and provide the prompts and questions, and/or transcript in the Supporting Information.**

# 1. Acknowledge

## Sample Acknowledgment Template

I acknowledge the use of **[insert AI system(s) and link]** to **[specific use of generative artificial intelligence].** The prompts used include **[list of prompts].** The output from these prompts was used to **[explain use].**

## 2. Remind

- Remind your coauthors, and yourself, that the output of the ChatGPT model is merely a very early draft, at best.

- The output is incomplete, might contain incorrect information, and every sentence and statement must be considered critically.

- Check, check, and check again.

- And then check again.

## 3. Do not use text verbatim

- Do not use text verbatim from ChatGPT.

- These are not your words. The bot might have also reused text from other sources, leading to inadvertent **plagiarism**.

## Sample APA Citation

**In the body of the essay:**

When prompted for a list of "public-service jobs for men," ChatGPT listed two stereotypically masculine jobs first: "Firefighter" and "Police Officer." (OpenAI, 2024).

**In the reference list:**

OpenAI. (2024). ChatGPT (Mar 14 version) [Large language model]. https://chat.openai.com/chat

# Sample MLA Citation

**In-text citation of quoted text where the prompt is described in the sentence:**

When prompted to list public-service jobs for men, ChatGPT listed two stereotypically masculine jobs first: "Firefighter" and "Police Officer." (2024).

**In-text citation of paraphrased text where the prompt is referenced in parentheses:**

ChatGPT has been known to output stereotypical answers when asked to list jobs for a certain gender. ("Public-service jobs for men" 2024).

**In the Works Cited list:**

"Give me a list of ten public-service jobs for men" prompt. ChatGPT 4o, version unknown, OpenAI, 3 Jul. 2024, https://chatgpt.com/share/78c36969-a...b-0085332182e8.

# 4. Verify with the original literature/source

Any citations recommended by an AI bot/ChatGPT **need to be verified** with the original literature/sources since the bot is known to generate erroneous citations.

## 5. Do not include ChatGPT as a co-author.

- AI systems cannot be attributed as authors. Only humans can.

- Do not include ChatGPT or any other AI-based bot as a co-author.

- It cannot generate new ideas or compose a discussion based on new results, as that is our domain as humans.

- It is merely a tool, like many other programs, for helping with the formulation and writing of manuscripts.

- Please refer to ACS Nano author guidelines for more information

# 6. ChatGPT cannot be held accountable

- ChatGPT cannot be held accountable for any statement or ethical breach.

- As it stands, **all authors of a manuscript share this responsibility.**

# 7. Do not allow ChatGPT to squelch

- Most importantly, do not allow ChatGPT to squelch your creativity and deep thinking.

- Use it to expand your horizons, and spark new ideas!

# How much AI content is acceptable in research/academic paper?

# Acceptable in a research paper.

| Field | Percentage of AI Usage |
|---|---|
| Natural Sciences (Physics, Chemistry, Biology) | 20-30% |
| Social Sciences (Psychology, Sociology, Anthropology) | 10-20% |
| Medical Sciences (Medicine, Nursing, Pharmacy) | 25-35% |
| Engineering (Computer Science, Electrical Engineering, Mechanical Engineering) | 30-40% |
| Humanities (English, History, Philosophy) | 5-10% |
| Education (Education Policy, Curriculum Studies, Educational Psychology) | 10-15% |
| Business (Marketing, Finance, Management) | 20-25% |
| Law (Legal Studies, Criminology, Criminal Justice) | 15-20% |
| Arts and Design (Fine Arts, Music, Architecture) | 5-10% |

# Position 3:
## Cloud Architect

With the rapid adoption of cloud computing, organizations are looking for skilled cloud architects who can design and manage cloud infrastructure. Cloud architects ensure the efficient utilization of cloud resources, implement security measures, and optimize performance.

Cagayan de Oro
ICT Business Council

# Position 1:
## Data Scientist

Data scientists are expected to be in high demand in the coming years. With the increasing reliance on data-driven decision making, companies are seeking professionals who can analyze large datasets, build machine learning models, and derive actionable insights.

Cagayan de Oro
ICT Business Council

# RECOMMENDATIONS

Have students work on projects that can serve as their portfolio when they apply for jobs. With IT students, constantly introduce new tools (e.g. project management) to them. For example, New IT graduates find it hard to compete with other applicants because they don't have anything to show how familiar they are of a specific programming language or tool.

END