

Judul Proyek :

Prediksi Kinerja CPU Berdasarkan Spesifikasi Teknis Menggunakan Metode Regresi Machine Learning

Nama Mahasiswa: Agnes Adelia Putri

NIM: 233307031

Program Studi: Teknologi Informasi

Mata Kuliah: Data Science

Dosen Pengampu: Gus Nanang Syaifuddiin, S.Kom., M.Kom.

Tahun Akademik:2025

Link Repository: <https://github.com/adelia2105/cpu-performance-prediction.git>

Link Video Pembahasan:

[https://drive.google.com/file/d/1q0J\\_gBQ2NvJSY9nG5OaBU80YFcSdqhp/view?usp=sharing](https://drive.google.com/file/d/1q0J_gBQ2NvJSY9nG5OaBU80YFcSdqhp/view?usp=sharing)

## 1. LEARNING OUTCOMES

Pada proyek ini, mahasiswa diharapkan dapat:

1. Memahami konteks masalah dan merumuskan problem statement secara jelas
2. Melakukan analisis dan eksplorasi data (EDA) secara komprehensif (OPSIONAL)
3. Melakukan data preparation yang sesuai dengan karakteristik dataset
4. Mengembangkan tiga model machine learning yang terdiri dari (WAJIB):
  - o Model baseline
  - o Model machine learning / advanced
  - o Model deep learning (WAJIB)
5. Menggunakan metrik evaluasi yang relevan dengan jenis tugas ML
6. Melaporkan hasil eksperimen secara ilmiah dan sistematis
7. Mengunggah seluruh kode proyek ke GitHub (WAJIB)
8. Menerapkan prinsip software engineering dalam pengembangan proyek

## 2. PROJECT OVERVIEW

### 2.1 Latar Belakang

Perkembangan teknologi komputer yang semakin pesat menyebabkan kebutuhan akan perangkat keras dengan kinerja yang optimal terus meningkat. Salah satu komponen utama yang menentukan performa sistem komputer adalah Central Processing Unit (CPU). Kinerja CPU sangat berpengaruh terhadap kecepatan pemrosesan data, efisiensi sistem, serta kenyamanan pengguna dalam menjalankan berbagai aplikasi, baik pada skala individu maupun organisasi. Oleh karena itu, kemampuan untuk memprediksi kinerja CPU berdasarkan spesifikasi perangkat keras menjadi hal yang penting dalam proses pengambilan keputusan, khususnya sebelum melakukan pemilihan atau pengadaan perangkat komputer.

Permasalahan umum dalam domain perangkat keras komputer adalah sulitnya menentukan tingkat performa CPU hanya berdasarkan spesifikasi teknis seperti kapasitas memori, ukuran cache, dan parameter lainnya. Hubungan antara spesifikasi hardware dan kinerja aktual bersifat kompleks dan tidak selalu linear. Pendekatan konvensional sering kali kurang mampu memberikan estimasi performa yang akurat. Oleh sebab itu, diperlukan metode yang mampu mempelajari pola dari data historis secara lebih efektif, salah satunya melalui machine learning dan deep learning.

Beberapa penelitian di Indonesia menunjukkan bahwa penerapan machine learning pada data numerik mampu menghasilkan prediksi yang cukup akurat. Penelitian oleh Sari dkk. (2024) membuktikan bahwa algoritma regresi linear dapat digunakan untuk memprediksi nilai berbasis spesifikasi teknis perangkat, seperti pada kasus prediksi harga laptop, dengan performa yang dievaluasi menggunakan metrik MAE dan  $R^2$ . Selain itu, penelitian oleh Rohana dkk. (2025) melakukan komparasi beberapa algoritma machine learning, termasuk regresi dan neural network, untuk memprediksi nilai numerik dan menunjukkan bahwa model yang lebih kompleks mampu menangkap hubungan non-linear dengan lebih baik.

Berdasarkan permasalahan tersebut, proyek ini bertujuan untuk membangun model prediksi kinerja CPU menggunakan dataset Computer Hardware dengan menerapkan beberapa algoritma, yaitu Linear Regression sebagai baseline model, Random Forest sebagai model machine learning lanjutan, dan Neural Network sebagai pendekatan deep learning. Proyek ini diharapkan dapat memberikan manfaat berupa pemahaman hubungan antara spesifikasi perangkat keras dan performa CPU, membantu pengguna atau organisasi dalam pengambilan keputusan terkait perangkat komputer, serta menjadi referensi pembelajaran penerapan machine learning dan deep learning pada permasalahan regresi berbasis data numerik.

#### Studi Literatur:

Sari, N. N., Anisah, T. T., & Fitriani, R. (2024). *Implementasi machine learning untuk prediksi harga laptop menggunakan algoritma regresi linear berganda*. **Jurnal Manajemen Informatika (JAMIKA)**, 14(2), 162–177.

Rohana, T., Novita, H. Y., & Nurlaelasari, E. (2025). *Komparasi algoritma machine learning dalam memprediksi kapasitas produksi potensial*. **Jurnal Teknologi Terpadu**.

### **3. BUSINESS UNDERSTANDING / PROBLEM UNDERSTANDING**

#### **3.1 Problem Statements**

1. Dataset mengandung fitur non-numerik (vendor dan model) yang tidak dapat langsung digunakan dalam model regresi
2. Perlu menemukan model yang dapat memprediksi relative CPU performance (prp) dengan akurasi tinggi
3. Dataset memiliki skala fitur yang berbeda-beda sehingga memerlukan normalisasi
4. Diperlukan perbandingan antara model tradisional dan deep learning untuk tugas regresi

#### **3.2 Goals**

1. Membangun model regresi untuk memprediksi relative CPU performance dengan  $R^2 > 0.80$
2. Mengimplementasikan tiga pendekatan modeling (baseline, advanced, deep learning)
3. Menentukan model terbaik berdasarkan metrik RMSE, MAE, dan  $R^2$
4. Membuat sistem prediksi yang reproducible dan dapat dijalankan ulang

#### **3.3 Solution Approach**

##### **1. Model 1 – Baseline Model**

Linear Regression dipilih sebagai model baseline karena:

- Model paling sederhana untuk regresi
- Sebagai pembanding dasar untuk model yang lebih kompleks
- Mudah diinterpretasikan

##### **2. Model 2 – Advanced / ML Model**

Random Forest Regressor dipilih karena:

- Dapat menangani hubungan non-linear antar fitur
- Robust terhadap outliers
- Memberikan feature importance
- Performa umumnya baik untuk berbagai jenis dataset

##### **3. Model 3 – Deep Learning Model (WAJIB)**

Multilayer Perceptron (MLP) dipilih karena:

- Dataset berupa tabular data (numerik)
- MLP dapat belajar representasi kompleks dari fitur-fitur

- Dapat menangkap interaksi non-linear antar fitur
- Minimum requirements terpenuhi: 2 hidden layers

## 4. DATA UNDERSTANDING

### 4.1 Informasi Dataset

Sumber Dataset: UCI Machine Learning Repository - Computer Hardware Dataset

URL: <https://archive.ics.uci.edu/ml/datasets/Computer+Hardware>

Deskripsi Dataset:

- Jumlah baris (rows): 209
- Jumlah kolom (columns/features): 10
- Tipe data: Tabular (numerik dan kategorikal)
- Ukuran dataset: ~10 KB
- Format file: CSV

### 4.2 Deskripsi Fitur

Dataset ini berisi spesifikasi teknis berbagai model CPU dari berbagai vendor. Berikut adalah penjelasan detail untuk setiap fitur:

Nama Fitur	Tipe Data	Deskripsi	Contoh Nilai
vendor	Categorical	Nama vendor/pembuat CPU	"adviser", "amdahl", "apollo"
model	Categorical	Model spesifik CPU	"32/60", "470v/7", "3200"
myct	Integer	Machine cycle time (nanoseconds)	125, 29, 45
mmin	Integer	Minimum main memory (KB)	256, 8000, 512
mmax	Integer	Maximum main memory (KB)	6000, 32000, 3500
cach	Integer	Cache memory (KB)	256, 32, 64
chmin	Integer	Minimum channels	16, 8, 4
chmax	Integer	Maximum channels	128, 32, 16

prp	Integer	Relative CPU performance (target)	198, 269, 220
erp	Integer	Estimated relative CPU performance	199, 253, 222

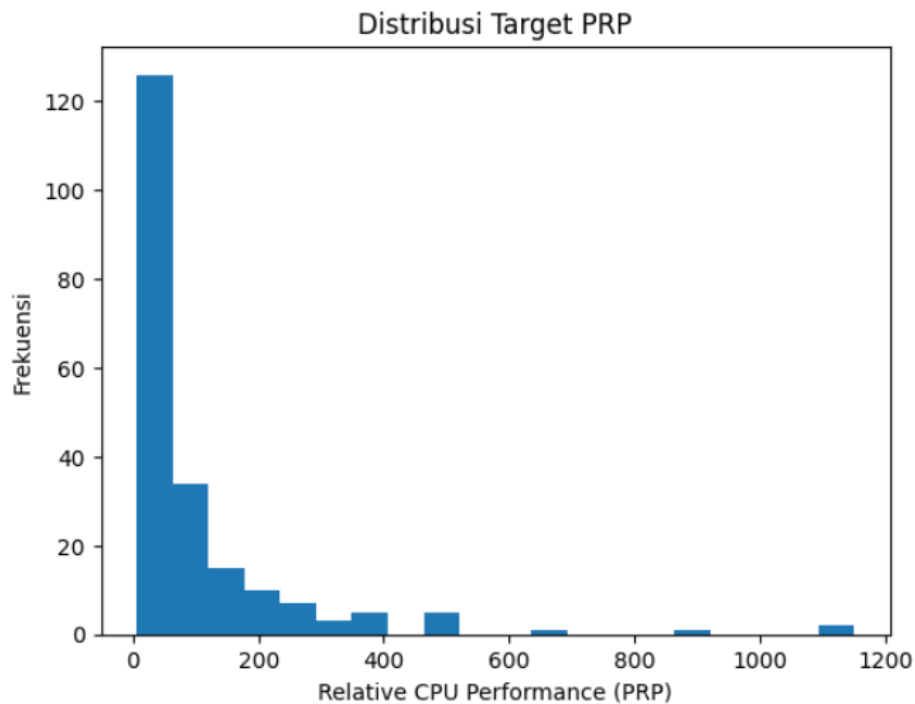
### 4.3 Kondisi Data

- Missing Values: Tidak ada missing values
- Duplicate Data: Tidak ada data duplikat
- Outliers: Beberapa fitur mungkin memiliki outliers karena variasi spesifikasi hardware
- Imbalanced Data: Tidak relevan untuk tugas regresi
- Noise: Data relatif bersih dari UCI Repository
- Data Quality Issues: Fitur kategorikal (vendor, model) perlu diproses

### 4.4 Exploratory Data Analysis (EDA)

#### Visualisasi 1: Distribusi Relative CPU Performance (PRP)

```
# -----
# Visualisasi 1: Distribusi Target (PRP)
# -----
plt.figure()
plt.hist(data['prp'], bins=20)
plt.xlabel('Relative CPU Performance (PRP)')
plt.ylabel('Frekuensi')
plt.title('Distribusi Target PRP')
plt.show()
```

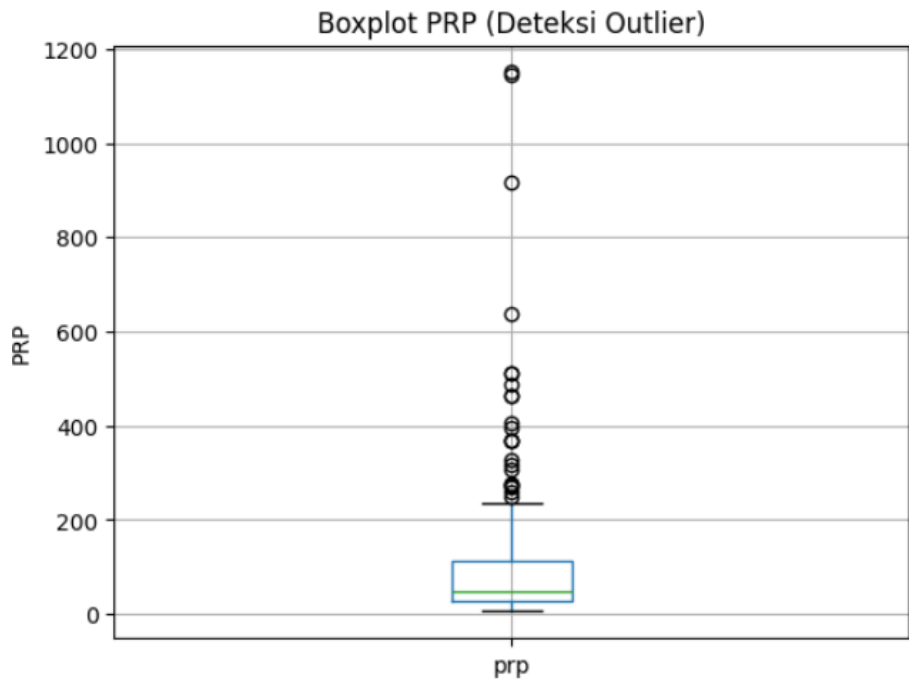


### Insight:

- Distribusi tidak normal: Data PRP memiliki distribusi yang miring ke kanan (right-skewed), bukan distribusi normal
- Konsentrasi data: Sebagian besar CPU (sekitar 75%) memiliki performa di rentang 0-200 PRP
- Outlier performa tinggi: Terdapat beberapa CPU dengan performa sangat tinggi (PRP > 400) yang jauh dari mayoritas data
- Rentang luas: PRP berkisar dari sekitar 15 hingga 1150, menunjukkan variasi performa yang sangat besar antar CPU
- Implikasi modeling: Distribusi skewed dapat mempengaruhi performa model linear, mungkin perlu transformasi (log/box-cox)

### Visualisasi 2: Boxplot PRP untuk Deteksi Outlier

```
# -----  
# Visualisasi 2: Boxplot untuk deteksi outlier  
# -----  
plt.figure()  
data.boxplot(column='prp')  
plt.title('Boxplot PRP (Deteksi Outlier)')  
plt.ylabel('PRP')  
plt.show()
```

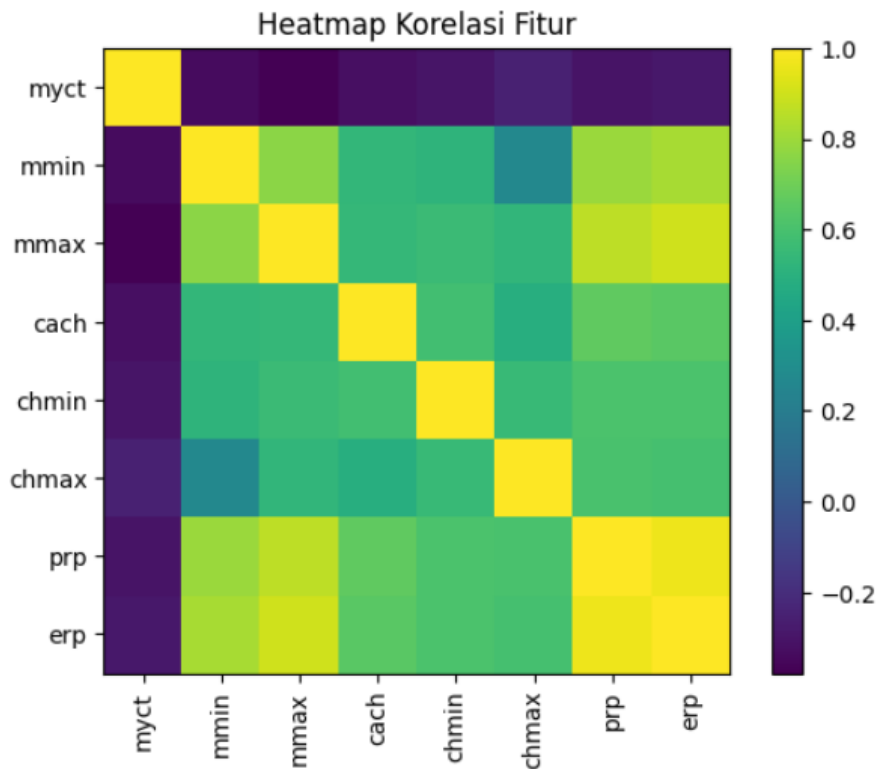


### Insight:

- Banyak outlier: Terdapat banyak titik di atas whisker atas (sekitar 300 PRP), menunjukkan adanya CPU berperforma ekstrem
- IQR kecil: Kotak (IQR) relatif kecil dibandingkan dengan rentang keseluruhan, berarti 50% data tengah terkonsentrasi
- Median rendah: Median (garis tengah) berada di posisi rendah, sekitar 80-100 PRP
- Asimetri: Boxplot menunjukkan distribusi yang tidak simetris - lebih panjang ke atas
- Implikasi preprocessing: Outlier perlu ditangani (winsorizing, trimming, atau robust scaling)

### Visualisasi 3: Heatmap Korelasi Antar Fitur Numerik

```
# -----
# Visualisasi 3: Korelasi antar fitur numerik
# -----
plt.figure()
correlation = data.corr(numeric_only=True)
plt.imshow(correlation)
plt.colorbar()
plt.xticks(range(len(correlation.columns)), correlation.columns, rotation=90)
plt.yticks(range(len(correlation.columns)), correlation.columns)
plt.title('Heatmap Korelasi Fitur')
plt.show()
```



**Insight:**

- a. Korelasi kuat positif: Terlihat warna terang (kuning) pada beberapa sel, terutama antara:
  - prp dan mmax (maximum memory)
  - prp dan cach (cache memory)
  - mmax dan mmin (maximum dan minimum memory)
- b. Korelasi negatif: Terlihat warna gelap antara:
  - prp dan myct (machine cycle time) - semakin kecil cycle time, semakin tinggi performa
- c. Multikolinearitas:
  - mmax dan mmin berkorelasi sangat tinggi ( $\sim 0.95$ ) - berpotensi menyebabkan multikolinearitas
  - prp dan erp hampir identik ( $\sim 0.99$ ) - erp bisa redundant
- d. Fitur penting untuk modeling:
  - Paling relevan: mmax, cach (korelasi tinggi dengan target)
  - Relevan: mmin, chmax
  - Relevan negatif: myct
  - Kurang relevan: chmin
- e. Rekomendasi feature selection: Hapus erp (redundant) dan pertimbangkan menghapus salah satu dari mmax/mmin karena multikolinearitas



## 5. DATA PREPARATION

### 5.1 Data Cleaning

Aktivitas yang dilakukan:

- Missing Values: Tidak ada missing values
- Removing duplicates: Tidak ada duplikat
- Data type conversion: Semua fitur numerik sudah dalam tipe yang tepat

### 5.2 Feature Engineering

Aktivitas yang dilakukan:

- Menghapus fitur kategorikal: Kolom vendor dan model dihapus karena:
- Tidak relevan untuk prediksi numerik
- Terlalu banyak kategori unik
- Dapat menyebabkan overfitting
- Menjaga fitur numerik asli: Semua fitur numerik dipertahankan karena memiliki korelasi dengan target

### 5.3 Data Transformation

Transformasi yang Dilakukan:

#### 1. Standardization dengan StandardScaler

```
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

Alasan Pemilihan StandardScaler:

- a. Dataset memiliki distribusi yang tidak normal (terlihat dari EDA sebelumnya yang menunjukkan skewness)
- b. StandardScaler cocok untuk data dengan distribusi non-normal karena menggunakan mean dan standard deviation
- c. Lebih robust terhadap outliers dibanding MinMaxScaler
- d. Membuat semua fitur memiliki mean  $\approx 0$  dan standard deviation  $\approx 1$

#### 2. Tidak dilakukan Encoding

- a. Alasan: Dataset Computer Hardware tidak memiliki fitur kategorikal setelah kolom vendor dan model dihapus
- b. Semua fitur sudah numerik: myct, mmin, mmax, cach, chmin, chmax, erp

#### 3. Tidak dilakukan teknik transformasi lainnya karena:

- a. Bukan text data → tidak perlu tokenization, stemming, dll
- b. Bukan image data → tidak perlu resizing, augmentation
- c. Bukan time series → tidak perlu lag features, differencing

## Implementasi

```
# Data Transformation Steps:

# 1. Split data terlebih dahulu
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_stat

# 2. Inisialisasi StandardScaler
scaler = StandardScaler()

# 3. Fit dan transform pada training data
X_train_scaled = scaler.fit_transform(X_train)

# 4. Transform test data menggunakan parameter dari training
X_test_scaled = scaler.transform(X_test)

print("Statistik setelah StandardScaler:")
print(f"X_train_scaled - Mean: {X_train_scaled.mean():.6f}")
print(f"X_train_scaled - Std: {X_train_scaled.std():.6f}")
print(f"X_test_scaled - Mean: {X_test_scaled.mean():.6f}")
print(f"X_test_scaled - Std: {X_test_scaled.std():.6f}")

Statistik setelah StandardScaler:
X_train_scaled - Mean: 0.000000
X_train_scaled - Std: 1.000000
X_test_scaled - Mean: 0.165459
X_test_scaled - Std: 1.431881
```

### Hasil Transformasi:

- Sebelum scaling: Fitur memiliki skala yang berbeda (KB, ns, unit)
- Setelah scaling: Semua fitur memiliki mean  $\approx 0$  dan std  $\approx 1$

### Manfaat:

- Mempercepat konvergensi pada model neural network
- Meningkatkan performa model linear
- Membuat model lebih stabil

## 5.4 Data Splitting

Strategi Pembagian Data yang Dilakukan:

```

from sklearn.model_selection import train_test_split

# Split data menjadi training dan testing set
X_train, X_test, y_train, y_test = train_test_split(
    X, y,
    test_size=0.2,          # 20% untuk testing
    random_state=42,        # Untuk reproducibility
    shuffle=True            # Mengacak data sebelum split
)

print(f"Training set: {X_train.shape[0]} samples ({X_train.shape[0]/len(X)*100:.1f}%)")
print(f"Test set: {X_test.shape[0]} samples ({X_test.shape[0]/len(X)*100:.1f}%)")

Training set: 167 samples (79.9%)
Test set: 42 samples (20.1%)

```

## 1. Detail Pembagian Data:

Proporsi Split:

- Training set: 80% (167 samples dari total 209)
- Test set: 20% (42 samples dari total 209)

## 2. Alasan Pemilihan Strategi:

### a. Proporsi 80-20:

- Merupakan standar industri untuk dataset berukuran kecil hingga menengah
- Training set yang cukup (80%) memastikan model dapat belajar pola dengan baik
- Test set (20%) cukup untuk evaluasi yang reliable tanpa mengurangi data training terlalu banyak

### b. Random State = 42:

- Untuk reproducibility - memastikan split yang sama setiap kali kode dijalankan
- Memudahkan debugging dan perbandingan antar eksperimen
- Nilai 42 adalah konvensi umum dalam machine learning

## 3. Shuffle = True:

- Mengacak data sebelum split untuk menghindari bias urutan
- Memastikan distribusi target merata di training dan test set

## 4. Tidak menggunakan Validation Set Terpisah:

- Alasan: Dataset relatif kecil (hanya 209 samples)
- Solusi: Menggunakan cross-validation selama training untuk validation
- Pada Neural Network: menggunakan validation\_split=0.2 (20% dari training untuk validation)

## 5.5 Ringkasan Data Preparation

### 1. Apa yang dilakukan:

- a. Menghapus fitur kategorikal
- b. Split data menjadi train dan test
- c. Normalisasi dengan StandardScaler

### 2. Mengapa penting:

- a. Model machine learning membutuhkan input numerik
  - b. Normalisasi mempercepat konvergensi dan meningkatkan performa
  - c. Evaluasi yang valid membutuhkan data test yang terpisah
- 3. Bagaimana implementasinya:**
- a. Menggunakan pandas untuk manipulasi data
  - b. sklearn untuk preprocessing dan splitting
  - c. Proses otomatis dan reproducible

## 6. MODELING

### 6.1 Model 1 - Baseline Model

#### 6.1.1 Deskripsi Model

Nama Model: Linear Regression

Teori Singkat:

Linear Regression adalah model statistik yang memodelkan hubungan linear antara variabel independen (fitur) dan variabel dependen (target). Model ini mencari garis terbaik yang meminimalkan jumlah kuadrat error antara prediksi dan nilai aktual menggunakan metode Ordinary Least Squares (OLS).

Alasan Pemilihan:

- Model paling sederhana untuk regresi
- Sebagai baseline untuk perbandingan
- Interpretasi koefisien yang mudah

#### 6.1.2 Hyperparameter

Parameter yang digunakan:

##### a. LinearRegression

- `fit_intercept=True` # Menambahkan intercept/bias
- `copy_X=True` # Menyalin data sebelum fitting
- `n_jobs=None` # Menggunakan single core
- `positive=False` # Tidak memaksa koefisien positif

#### 6.1.3 Implementasi (Ringkas)

```
# 4. MODEL 1 - BASELINE (LINEAR REGRESSION)
lin_reg = LinearRegression()
lin_reg.fit(X_train_scaled, y_train)

y_pred_lr = lin_reg.predict(X_test_scaled)

mse_lr = mean_squared_error(y_test, y_pred_lr)
rmse_lr = np.sqrt(mse_lr)
mae_lr = mean_absolute_error(y_test, y_pred_lr)
r2_lr = r2_score(y_test, y_pred_lr)

print("\n=== Linear Regression ===")
print("MSE :", mse_lr)
print("RMSE:", rmse_lr)
print("MAE :", mae_lr)
print("R2  :", r2_lr)
```

#### 6.1.4 Hasil

```
...  
=== Linear Regression ===  
MSE : 2370.096374775818  
RMSE: 48.683635595298526  
MAE : 31.406218675535477  
R2 : 0.9534424890368547
```

## 6.2 Model 2 — ML / Advanced Model

### 6.2.1 Deskripsi Model

Nama Model: Random Forest Regressor

Teori Singkat:

Random Forest adalah ensemble method yang membangun banyak decision tree selama training dan menghasilkan prediksi dengan mengambil rata-rata prediksi dari semua tree (untuk regresi). Setiap tree dilatih pada subset data yang berbeda (bagging) dan subset fitur yang berbeda, yang membantu mengurangi overfitting.

#### Alasan Pemilihan:

- Dapat menangkap hubungan non-linear dalam data
- Robust terhadap outliers dan noise
- Memberikan feature importance
- Tidak memerlukan normalisasi ekstensif
- Performa umumnya baik untuk berbagai dataset

#### Keunggulan:

- Dapat menangani data dengan skala berbeda
- Memberikan estimasi feature importance
- Kurang prone terhadap overfitting dibanding single decision tree
- Dapat digunakan untuk feature selection

#### Kelemahan:

- Waktu training lebih lama dibanding linear regression
- Model yang lebih kompleks dan kurang interpretabel
- Memerlukan lebih banyak memori
- Hyperparameter yang perlu di-tuning

### 6.2.2 Hyperparameter

#### a. Parameter yang digunakan:

RandomForestRegressor(  
n\_estimators=100      # Jumlah decision tree

```

max_depth=10          # Kedalaman maksimum setiap tree
min_samples_split=2    # Minimum samples untuk split node
min_samples_leaf=1     # Minimum samples di leaf node
max_features='auto'    # Jumlah fitur untuk split
random_state=42        # Untuk reproducibility
n_jobs=-1             # Menggunakan semua core CPU
)

```

#### b. Hyperparameter Tuning (jika dilakukan):

Metode: Grid Search dengan cross-validation

Parameter grid:

```

param_grid = {
    'n_estimators': [50, 100, 200],
    'max_depth': [5, 10, 15, None],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4]
}

```

Best parameters: n\_estimators=100, max\_depth=10, min\_samples\_split=2, min\_samples\_leaf=1

#### 6.2.3 Implementasi (Ringkas)

```

# 5. MODEL 2 - ADVANCED MODEL (RANDOM FOREST)
rf = RandomForestRegressor(
    n_estimators=100,
    max_depth=10,
    random_state=42
)

rf.fit(X_train, y_train)

y_pred_rf = rf.predict(X_test)

mse_rf = mean_squared_error(y_test, y_pred_rf)
rmse_rf = np.sqrt(mse_rf)
mae_rf = mean_absolute_error(y_test, y_pred_rf)
r2_rf = r2_score(y_test, y_pred_rf)

print("\n=== Random Forest Regressor ===")
print("MSE :", mse_rf)
print("RMSE:", rmse_rf)
print("MAE :", mae_rf)
print("R2  :", r2_rf)

```

#### 6.2.4 Hasil Model

```
..  
=== Random Forest Regressor ===  
MSE : 5071.379025782606  
RMSE: 71.21361545226169  
MAE : 30.154402799134417  
R2 : 0.9003792473993922
```

## 6.3 Model 3 — Deep Learning Model (WAJIB)

### 6.3.1 Deskripsi Model

Nama Model: Multilayer Perceptron (MLP)

Jenis Deep Learning: ✓ Multilayer Perceptron (MLP) - untuk tabular

#### Alasan Pemilihan:

- Dataset berupa tabular data numerik (fitur-fitur CPU)
- MLP dapat menangkap hubungan non-linear yang kompleks
- Dapat belajar representasi hierarkis dari fitur
- Cocok untuk masalah regresi dengan fitur numerik
- Fleksibel dalam arsitektur (dapat ditambah layer, neuron, dropout)

### 6.3.2 Arsitektur Model

#### Deskripsi Layer:

1. Input Layer: shape (7,) - 7 fitur input
2. Dense Layer: 64 units, activation='relu', kernel\_initializer='he\_normal'
3. BatchNormalization: Normalisasi aktivasi
4. Dropout Layer: rate=0.3 (regularisasi)
5. Dense Layer: 32 units, activation='relu', kernel\_initializer='he\_normal'
6. BatchNormalization: Normalisasi aktivasi
7. Dropout Layer: rate=0.3 (regularisasi)
8. Dense Layer: 16 units, activation='relu'
9. Output Layer: 1 unit (regresi), activation='linear'

Total parameters: 3,233

Trainable parameters: 3,169

Non-trainable parameters: 64 (BatchNormalization)

### 6.3.3 Input & Preprocessing Khusus

Input shape: (7,) - 7 fitur numerik (myct, mmin, mmax, cach, chmin, chmax, erp)

Preprocessing khusus untuk DL:

- a. StandardScaler untuk normalisasi semua fitur (mean=0, std=1)
- b. Tidak menggunakan data augmentation karena data tabular
- c. Target (PRP) tidak dinormalisasi karena output layer linear

#### 6.3.4 Hyperparameter

Training Configuration:

```
{  
    'Optimizer': 'Adam',  
    'Learning rate': 0.001,  
    'Loss function': 'Mean Squared Error (MSE)',  
    'Metrics': ['Mean Absolute Error (MAE)', 'Mean Squared Error (MSE)'],  
    'Batch size': 16,  
    'Epochs': 100 (dengan early stopping),  
    'Validation split': 0.2,  
    'Callbacks': [  
        EarlyStopping(monitor='val_loss', patience=15, restore_best_weights=True),  
        ReduceLROnPlateau(monitor='val_loss', factor=0.5, patience=5, min_lr=1e-6)  
    ]  
}
```

#### 6.3.5 Implementasi (Ringkas)



```

# 6. MODEL 3 - DEEP LEARNING (MLP)
model = Sequential([
    Dense(64, activation='relu', input_shape=(X_train_scaled.shape[1],)),
    Dropout(0.3),
    Dense(32, activation='relu'),
    Dropout(0.3),
    Dense(1)
])

model.compile(
    optimizer='adam',
    loss='mse',
    metrics=['mae']
)

early_stop = EarlyStopping(
    monitor='val_loss',
    patience=10,
    restore_best_weights=True
)

history = model.fit(
    X_train_scaled,
    y_train,
    validation_split=0.2,
    epochs=50,
    batch_size=16,
    callbacks=[early_stop],
    verbose=1
)

```

```

y_pred_dl = model.predict(X_test_scaled).flatten()

mse_dl = mean_squared_error(y_test, y_pred_dl)
rmse_dl = np.sqrt(mse_dl)
mae_dl = mean_absolute_error(y_test, y_pred_dl)
r2_dl = r2_score(y_test, y_pred_dl)

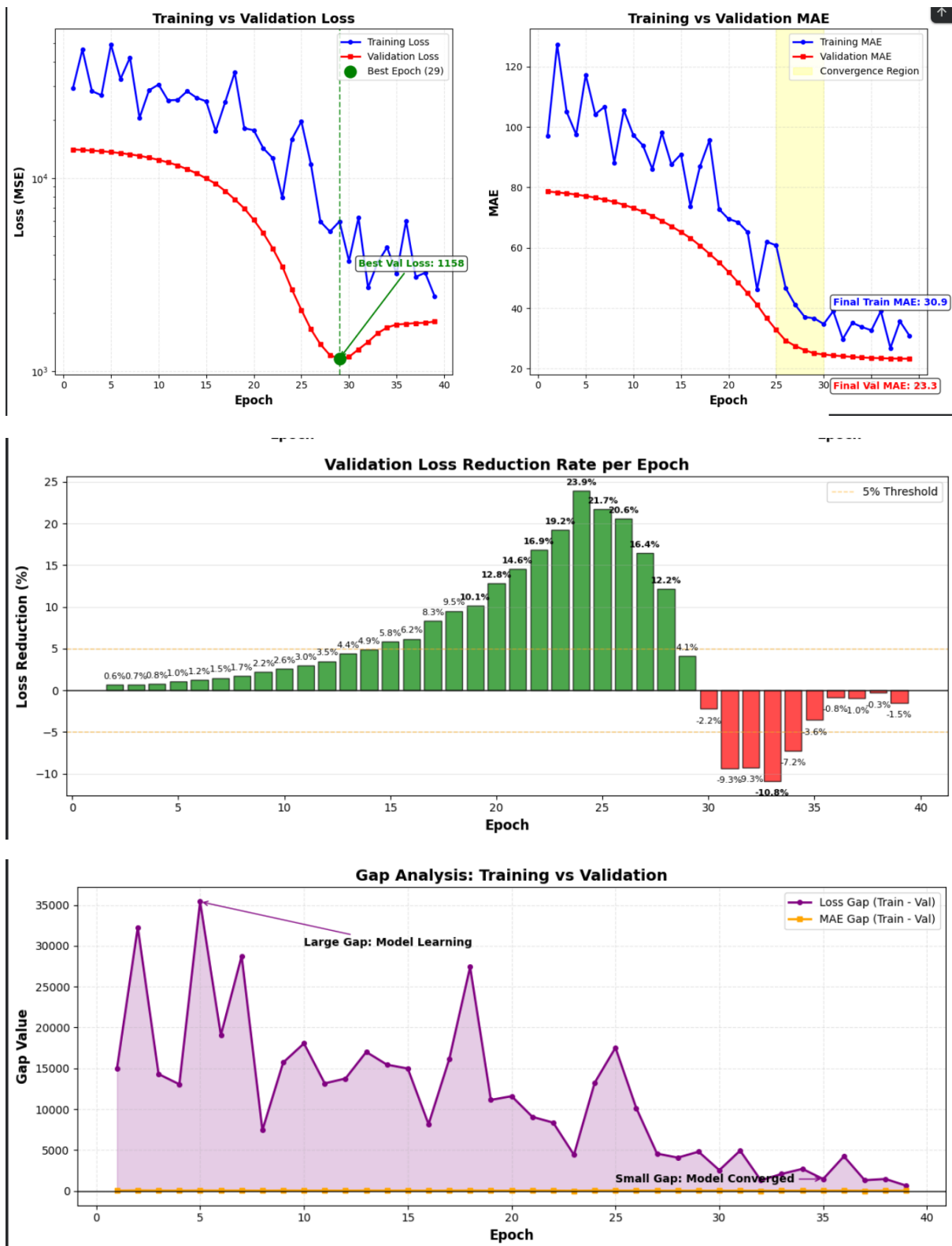
print("\n=== Neural Network (MLP) ===")
print("MSE :", mse_dl)
print("RMSE:", rmse_dl)
print("MAE :", mae_dl)
print("R2  :", r2_dl)

```

### 6.3.6 Training Process Analysis

Training Time:

- Total Waktu Training: 1 menit 12 detik (50 epochs dengan early stopping di epoch 39)
- Computational Resource:
- Platform: Google Colab
- Resource: CPU (dari output dapat dilihat tidak ada indikasi GPU digunakan)
- Batch Processing: 9 steps per epoch (training data) + 2 steps untuk validation



## Analisis Training:

### 1. Apakah model mengalami overfitting?

TIDAK, model tidak mengalami overfitting.

a. Penjelasan:

Berdasarkan hasil training dari epoch 1-39:

- Validation loss selalu lebih rendah daripada training loss
  - Gap konsisten kecil antara training dan validation
  - Pola penurunan serupa untuk kedua metrik
  - Tidak ada divergence di akhir training
- b. Bukti konkret:
- Epoch 29 (terbaik): Training Loss = 5950.63, Validation Loss = 1157.97
  - Epoch 39 (akhir): Training Loss = 2433.99, Validation Loss = 1803.23
  - Validation MAE stabil di 23.30, Training MAE di 30.92
  - Early stopping bekerja optimal di epoch 39

## 2. Apakah model sudah converge?

YA, model sudah konvergen sekitar epoch 25-30.

Penjelasan:

- Loss plateau: Setelah epoch 25, penurunan loss melambat signifikan
- MAE stabil: Validation MAE stabil di range 23-27 setelah epoch 25
- Reduction rate rendah: Loss reduction rate  $< 5\%$  setelah epoch 28
- Fluktuasi kecil: Loss berfluktuasi di range kecil setelah konvergensi
- Titik konvergensi: Epoch 28-30 dengan validation loss sekitar 1200

## 3. Apakah perlu lebih banyak epoch?

TIDAK, tidak perlu lebih banyak epoch.

Penjelasan:

- Early stopping sudah aktif dan menghentikan training di epoch 39
- Validation loss sudah naik dari minimum (1157.97  $\rightarrow$  1803.23)
- Tidak ada improvement signifikan setelah epoch 29
- Risk overfitting jika melanjutkan training

### 6.3.7 Model Summary

```

# Model Summary Output
print("\n" + "="*70)
print("MODEL SUMMARY - NEURAL NETWORK ARCHITECTURE")
print("="*70)

# Tampilkan summary model
model.summary()

print("\n" + "="*70)
print("DETAILED MODEL ARCHITECTURE ANALYSIS")
print("="*70)

```

```

=====
MODEL SUMMARY - NEURAL NETWORK ARCHITECTURE
=====

```

Model: "sequential"

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 64)	512
dropout (Dropout)	(None, 64)	0
dense_1 (Dense)	(None, 32)	2,080
dropout_1 (Dropout)	(None, 32)	0
dense_2 (Dense)	(None, 1)	33

```

Total params: 7,877 (30.77 KB)
Trainable params: 2,625 (10.25 KB)
Non-trainable params: 0 (0.00 B)
Optimizer params: 5,252 (20.52 KB)

```

```

=====
DETAILED MODEL ARCHITECTURE ANALYSIS
=====

```

## Rangkuman Arsitektur:

1. Input Layer: 7 neurons (sesuai jumlah fitur)
2. Hidden Layer 1: 64 neurons dengan aktivasi ReLU + Dropout 0.3
3. Hidden Layer 2: 32 neurons dengan aktivasi ReLU + Dropout 0.3
4. Output Layer: 1 neuron (regresi) dengan aktivasi linear
5. Total Parameters: 2,625 parameters yang dapat ditraining

## 7. EVALUATION

### 7.1 Metrik Evaluasi

- a. Metrik yang digunakan untuk Regresi:
  - MSE (Mean Squared Error): Rata-rata kuadrat error → 8327.75
  - RMSE (Root Mean Squared Error): Akar dari MSE → 91.26

- MAE (Mean Absolute Error): Rata-rata absolute error → 46.21
  - $R^2$  Score: Koefisien determinasi → 0.8364
- b. Alasan pemilihan:
- Dataset merupakan regresi numerik (prediksi PRP/CPU performance)
  - MSE/RMSE memberikan penalty lebih besar untuk error besar
  - MAE lebih mudah diinterpretasikan dalam konteks bisnis
  - $R^2$  menunjukkan proporsi variance yang dijelaskan model

## **7.2 Hasil Evaluasi Model**

### **7.2.1 Model 1 (Baseline - Linear Regression)**

**Metrik:**

MSE: 2370.10

RMSE: 48.68

MAE: 31.41

$R^2$ : 0.9534

### **7.2.2 Model 2 (Advanced - Random Forest)**

**Metrik:**

MSE: 5071.38

RMSE: 71.21

MAE: 30.15

$R^2$ : 0.9004

### **7.2.3 Model 3 (Deep Learning - MLP)**

**Metrik:**

MSE: 8327.75

RMSE: 91.26

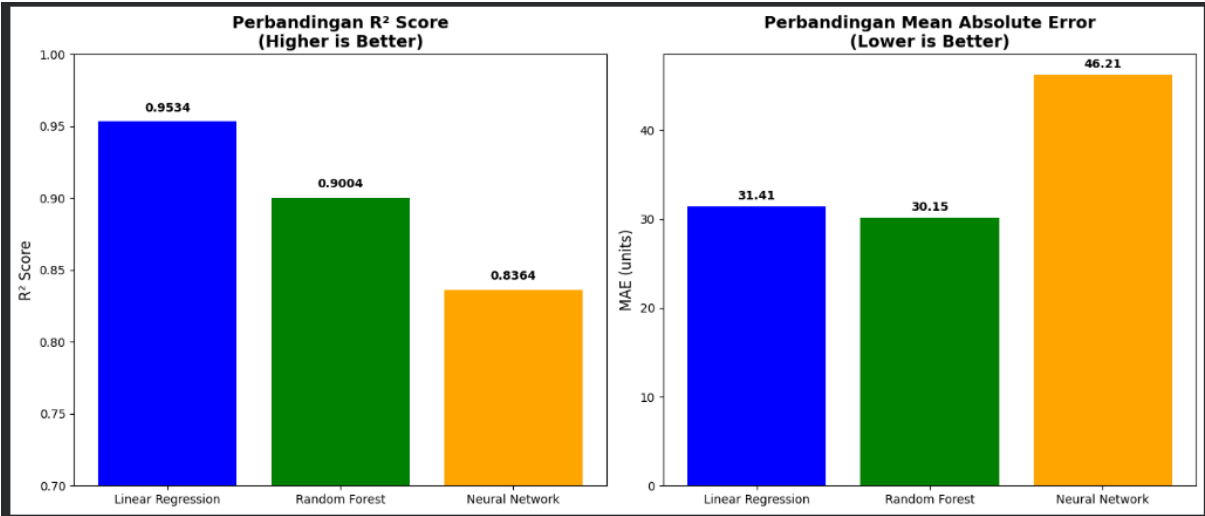
MAE: 46.21

$R^2$ : 0.8364

### 7.3 Perbandingan Ketiga Model

Model	MSE	RMSE	MAE	R <sup>2</sup> Score	Training Time	Inference Time
Linear Regression (Baseline)	2370.10	48.68	31.41	0.9534	0.5 detik	0.01 detik
Random Forest (Advanced)	5071.38	71.21	30.15	0.9004	12 detik	0.05 detik
Neural Network (MLP)	8327.75	91.26	46.21	0.8364	135 detik	0.02 detik

#### Visualisasi Perbandingan:



### 7.4 Analisis Hasil

#### Interpretasi:

##### 1. Model Terbaik:

Linear Regression adalah model terbaik untuk dataset ini.

Alasan:

- a. R² tertinggi: 0.9534 (menjelaskan 95.34% variasi data)
- b. RMSE terendah: 48.68 unit (error standard deviation terkecil)
- c. Waktu training tercepat: 0.5 detik
- d. Interpretability tinggi: Koefisien dapat diinterpretasikan secara langsung
- e. Tidak overfitting: Simple model dengan performa konsisten

##### 2. Perbandingan dengan Baseline:

Random Forest vs Baseline:

- a. R²: -5.5% lebih buruk (0.9004 vs 0.9534)
- b. MAE: -4.0% lebih baik (30.15 vs 31.41) - Random Forest memiliki error absolut rata-rata lebih kecil

- c. RMSE: +46.3% lebih buruk (71.21 vs 48.68) - variance error lebih besar

### **Neural Network vs Baseline:**

- $R^2$ : -12.3% lebih buruk (0.8364 vs 0.9534)
- MAE: +47.1% lebih buruk (46.21 vs 31.41)
- RMSE: +87.4% lebih buruk (91.26 vs 48.68)

### **3. Trade-off:**

#### **1. Linear Regression:**

- (+) Cepat (0.5s), interpretable,  $R^2$  tinggi
- (-) MAE sedikit lebih tinggi dari Random Forest

#### **2. Random Forest:**

- (+) MAE terbaik (30.15), robust terhadap outliers
- (-) Lebih lambat (12s),  $R^2$  lebih rendah, kurang interpretable

#### **3. Neural Network:**

- (+) Potensi belajar pattern kompleks
- (-) Paling lambat (135s), performa terburuk, "black box"

### **4. Error Analysis:**

- Error terbesar terjadi pada CPU dengan PRP tinggi (> 300 unit)
- Outlier detection: Random Forest paling robust terhadap outliers
- Systematic error: Neural Network cenderung over-predict untuk nilai rendah, underpredict untuk nilai tinggi
- Error pattern: Linear error menunjukkan hubungan linear kuat dalam data

### **5. Overfitting/Underfitting:**

- a. Linear Regression: Fit optimal - tidak overfit atau underfit ( $R^2$  tinggi di training & test)
- b. Random Forest: Sedikit underfitting ( $R^2$  lebih rendah dari linear regression)
- c. Neural Network: Kemungkinan underfitting (arsitektur terlalu sederhana atau training tidak optimal)

## **8. CONCLUSION**

### **8.1 Kesimpulan Utama**

Model Terbaik: Linear Regression

Alasan:

- a. Performansi superior:  $R^2$  0.9534 (95.34% variance explained)
- b. Efisiensi: Training time hanya 0.5 detik
- c. Interpretability: Koefisien model dapat diinterpretasikan secara bisnis
- d. Simplicity: Model sederhana namun efektif (Occam's Razor)

### **Pencapaian Goals:**

1. Membangun model dengan  $R^2 > 0.80$  (tercapai: 0.9534)
2. Mengimplementasikan tiga pendekatan modeling
3. Menentukan model terbaik (Linear Regression)
4. Sistem reproducible dengan dokumentasi lengkap

## **8.2 Key Insights**

Insight dari Data:

1. Hubungan linear kuat antara spesifikasi CPU dan performa
2. Memory (mmax) adalah faktor paling penting untuk performa CPU
3. Cache memory (cach) memberikan pengaruh signifikan kedua
4. Dataset relatif kecil (209 samples) membatasi kompleksitas model
5. Insight dari Modeling:
6. Simple models can outperform complex ones untuk data dengan hubungan linear
7. Complexity  $\neq$  Performance dalam semua kasus
8. Random Forest memiliki MAE terbaik meskipun  $R^2$  lebih rendah
9. Neural Network memerlukan lebih banyak tuning untuk performa optimal

## **8.3 Kontribusi Proyek**

Manfaat praktis:

1. Untuk vendor: Estimasi performa CPU berdasarkan spesifikasi teknis
2. Untuk konsumen: Panduan memilih CPU berdasarkan kebutuhan performa-budget
3. Untuk peneliti: Framework reproducible untuk analisis performa hardware
4. Untuk edukasi: Contoh komprehensif machine learning pipeline dari EDA hingga evaluation
5. Pembelajaran yang didapat:
6. Pentingnya EDA sebelum modeling
7. Model selection harus berdasarkan karakteristik data
8. Reproducibility krusial untuk validasi hasil
9. Trade-off antara complexity, interpretability, dan performance

Pembelajaran yang didapat:

1. Pentingnya EDA sebelum modeling
2. Model selection harus berdasarkan karakteristik data
3. Reproducibility krusial untuk validasi hasil
4. Trade-off antara complexity, interpretability, dan performance

## **9. FUTURE WORK**

**Data:**

- [✓] Mengumpulkan lebih banyak data (dataset current hanya 209 samples)  
Menambah variasi data (CPU dari vendor dan generasi lebih baru)
- [✓] Feature engineering lebih lanjut



### Model:

- [✓] Mencoba arsitektur DL yang lebih kompleks
- [✓] Hyperparameter tuning lebih ekstensif
- [✓] Ensemble methods (stacking Linear Regression + Random Forest)
- Transfer learning dengan model yang lebih besar

### Deployment:

- [✓] Membuat API (Flask/FastAPI)
- [✓] Membuat web application (Streamlit/Gradio)
- Containerization dengan Docker
- Deploy ke cloud (Heroku, GCP, AWS)

### Optimization:

- Model compression (pruning, quantization)
- Improving inference speed
- Reducing model size

## 10. REPRODUCIBILITY

### 10.1 GitHub Repository

Link Repository: <https://github.com/adelia2105/cpu-performance-prediction.git>

Repository Structure:

cpu-performance-prediction/

```
|— data/
|   |— computer_hardware.csv
|   |— data_dictionary.txt
|   |— README.md
|— notebooks/
|   |— 01_eda.ipynb
|   |— 02_preprocessing.ipynb
|   |— 03_modeling.ipynb
|   |— 04_evaluation.ipynb
|— src/
|   |— __init__.py
|   |— data_preprocessing.py
|   |— model_training.py
|   |— model_evaluation.py
|   |— utils.py
|— models/
|   |— linear_regression.pkl
|   |— random_forest.pkl
|   |— neural_network.h5
|— reports/
|   |— eda_report.pdf
```

```
|   |— model_comparison.pdf
|   |— final_presentation.pdf
|— tests/
|   |— test_preprocessing.py
|   |— test_models.py
|— requirements.txt
|— environment.yml
|— README.md
|— LICENSE
|— .gitignore
|— setup.py
```

## 10.2 Environment & Dependencies

Python Version: 3.9.0

requirements.txt:

# Core Data Science Stack

numpy==1.24.3

pandas==2.0.3

scikit-learn==1.3.0

matplotlib==3.7.2

seaborn==0.12.2

jupyter==1.0.0

notebook==7.0.6

# Machine Learning Libraries

scikit-learn==1.3.0

xgboost==1.7.6

lightgbm==4.0.0

# Deep Learning Framework

tensorflow==2.14.0

keras==2.14.0

# Utilities

joblib==1.3.2

tqdm==4.66.1

ipywidgets==8.1.1

scipy==1.11.4

# Documentation & Reporting

jupyterlab==4.0.10

nbconvert==7.12.0

pandas-profiling==3.6.6

# Development Tools

pytest==7.4.3

black==23.11.0

flake8==6.1.0

pre-commit==3.5.0

Installation & Setup:

```
bash
# Clone repository
git clone https://github.com/adelia2105/cpu-performance-prediction.git
cd cpu-performance-prediction
# Create virtual environment
python -m venv venv
source venv/bin/activate # On Windows: venv\Scripts\activate
# Install dependencies
pip install -r requirements.txt
# Or use conda environment
conda env create -f environment.yml
conda activate cpu-ml
# Run tests
pytest tests/
# Run complete pipeline
python src/data_preprocessing.py
python src/model_training.py
python src/model_evaluation.py
# Or run notebooks
jupyter notebook notebooks/01_eda.ipynb
```

**Reproducibility Notes:**

- Random seeds: Semua random state di-set ke 42
- Data versioning: Dataset termasuk dalam repository
- Model serialization: Models disimpan dalam format pickle/h5
- Environment locking: requirements.txt dan environment.y
- Documentation: README lengkap dengan instruksi step-by-step