# modeling_pipeline

November 23, 2025

# 1 Cloud Autoscaling – ML Forecasting Pipeline

## 1.1 Production-Grade Predictive Modeling Using Google Cluster 2019 Data

---

### 1.1.1 Overview

This notebook implements a **production-grade machine learning pipeline** for forecasting CPU demand using real Google Cluster 2019 trace data.

**Key Features:** - No data leakage (proper feature engineering with shifts) - Train/Validation/Test split (70/15/15) - Multiple models (Linear Regression, Random Forest, LightGBM) - Comprehensive evaluation and visualization - Multi-step recursive forecasting - Integration-ready outputs

**Objective:** Forecast CPU demand 1-12 steps ahead (5-60 minutes) to enable proactive autoscaling.

---

## 1.2 1. Imports and Setup

```python
import sys
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from pathlib import Path
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LinearRegression
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
from lightgbm import LGBMRegressor
import json
import warnings
warnings.filterwarnings('ignore')

# Add parent directory to path
sys.path.insert(0, str(Path.cwd().parent))
```

```python
# Import project loaders
from cloud_autoscale.data import GCP2019Loader

# Configure plotting
sns.set_style('whitegrid')
plt.rcParams['figure.figsize'] = (14, 6)
plt.rcParams['font.size'] = 11

print(' All imports successful')
print(f' Working directory: {Path.cwd()}')
```

```
 All imports successful
 Working directory: /Users/medhatabouzeid/Documents/00-Projects/_AUS/Cloud-
AutoScale/notebooks
```

## 1.3  2. Load GCP 2019 Data

Loading real Google Cluster 2019 traces (no synthetic data).

```python
[2]: # Load GCP 2019 data - FULL TRACE
     print('Loading GCP 2019 cluster trace data...')
     print('='*70)

     loader = GCP2019Loader(
         processed_dir='../data/processed',
         step_minutes=5,
         duration_minutes=None  # Use full trace
     )

     df = loader.load()

     print(f' Loaded {len(df):,} time steps')
     print(f' Time span: {df["time"].min():.0f} to {df["time"].max():.0f} minutes')
     print(f' Duration: {(df["time"].max() - df["time"].min()) / 60:.1f} hours')
     print(f'\nColumns: {list(df.columns)}')
     print('\nFirst 5 rows:')
     df.head()
```

```
Loading GCP 2019 cluster trace data…
======================================================================
 Loaded 8,929 time steps
 Time span: 0 to 44640 minutes
 Duration: 744.0 hours

Columns: ['step', 'time', 'cpu_demand', 'mem_demand', 'new_instances',
'new_instances_norm', 'machines_reporting']

First 5 rows:
```

```
[2]:    step  time  cpu_demand  mem_demand  new_instances  new_instances_norm  \
     0     0     0   10.283436    4.215649            0.0            0.000000
     1     1     5   11.116977    4.543966        14451.0            9.578588
     2     2    10   10.353116    4.374648        15688.0            9.660715
     3     3    15   12.320097    4.651689        13254.0            9.492130
     4     4    20   12.255638    4.914910        12053.0            9.397152

        machines_reporting
     0              2195.0
     1              2209.0
     2              2218.0
     3              2221.0
     4              2223.0
```

## 1.4   3. Feature Engineering (NO DATA LEAKAGE)

**Critical Fix:** All rolling windows are shifted BEFORE rolling to prevent data leakage.

### 1.4.1   Features:

1. **Lag Features** - Previous values (1, 2, 3, 6, 12 steps)
2. **Rolling Statistics** - Moving averages (SHIFTED first)
3. **Differencing** - Rate of change
4. **Cyclical** - Daily patterns

```python
[3]: # Create features with NO DATA LEAKAGE
     print('Creating features (preventing data leakage)...')
     print('='*70)

     df_features = df.copy()

     # 1. Lag Features
     print('\n[1/4] Lag features...')
     for lag in [1, 2, 3, 6, 12]:
         df_features[f'cpu_lag{lag}'] = df_features['cpu_demand'].shift(lag)
         df_features[f'mem_lag{lag}'] = df_features['mem_demand'].shift(lag)
         df_features[f'evt_lag{lag}'] = df_features['new_instances_norm'].shift(lag)
     print('   Created 15 lag features')

     # 2. Rolling Statistics (SHIFT FIRST to prevent leakage)
     print('\n[2/4] Rolling statistics (shifted to prevent leakage)...')
     for w in [3, 6, 12]:
         # CRITICAL: shift(1) BEFORE rolling to prevent data leakage
         df_features[f'cpu_ma{w}'] = df_features['cpu_demand'].shift(1).
      →rolling(window=w, min_periods=1).mean()
         df_features[f'mem_ma{w}'] = df_features['mem_demand'].shift(1).
      →rolling(window=w, min_periods=1).mean()
```

```
    df_features[f'evt_ma{w}'] = df_features['new_instances_norm'].shift(1).
 ↪rolling(window=w, min_periods=1).mean()
print('    Created 9 rolling features (no leakage)')

# 3. Differencing
print('\n[3/4] Differencing...')
df_features['cpu_diff1'] = df_features['cpu_demand'].diff()
df_features['mem_diff1'] = df_features['mem_demand'].diff()
print('    Created 2 differencing features')

# 4. Cyclical time features
print('\n[4/4] Cyclical features...')
df_features['sin_day'] = np.sin(2 * np.pi * df_features['step'] / 288)
df_features['cos_day'] = np.cos(2 * np.pi * df_features['step'] / 288)
print('    Created 2 cyclical features')

# Drop NaN
print('\n[5/5] Cleaning...')
rows_before = len(df_features)
df_clean = df_features.dropna().reset_index(drop=True)
rows_after = len(df_clean)

print(f'    Rows before: {rows_before:,}')
print(f'    Rows after: {rows_after:,}')
print(f'    Dropped: {rows_before - rows_after:,}')
print(f'\n Total features: {len(df_clean.columns)}')
print('='*70)
```

```
Creating features (preventing data leakage)…
======================================================================

[1/4] Lag features…
    Created 15 lag features

[2/4] Rolling statistics (shifted to prevent leakage)…
    Created 9 rolling features (no leakage)

[3/4] Differencing…
    Created 2 differencing features

[4/4] Cyclical features…
    Created 2 cyclical features

[5/5] Cleaning…
    Rows before: 8,929
    Rows after: 8,917
    Dropped: 12
```

```
    Total features: 35
    ======================================================================
```

## 1.5 4. Train/Validation/Test Split

**Split Strategy:** - Train: 70% (for model training) - Validation: 15% (for hyperparameter tuning) - Test: 15% (for final evaluation)

**Ordered split** to preserve temporal structure.

```python
[4]: # Define split indices
     total = len(df_clean)
     train_end = int(total * 0.7)
     val_end = int(total * 0.85)

     train = df_clean.iloc[:train_end]
     val = df_clean.iloc[train_end:val_end]
     test = df_clean.iloc[val_end:]

     print('='*70)
     print('TRAIN/VALIDATION/TEST SPLIT')
     print('='*70)
     print(f'\nTotal samples: {total:,}')
     print(f'\nTrain: {len(train):,} samples ({len(train)/total*100:.1f}%)')
     print(f'  Time: {train["time"].min():.0f} - {train["time"].max():.0f} min')
     print(f'\nValidation: {len(val):,} samples ({len(val)/total*100:.1f}%)')
     print(f'  Time: {val["time"].min():.0f} - {val["time"].max():.0f} min')
     print(f'\nTest: {len(test):,} samples ({len(test)/total*100:.1f}%)')
     print(f'  Time: {test["time"].min():.0f} - {test["time"].max():.0f} min')
     print('='*70)
```

```
    ======================================================================
    TRAIN/VALIDATION/TEST SPLIT
    ======================================================================

    Total samples: 8,917

    Train: 6,241 samples (70.0%)
      Time: 60 - 31260 min

    Validation: 1,338 samples (15.0%)
      Time: 31265 - 37950 min

    Test: 1,338 samples (15.0%)
      Time: 37955 - 44640 min
    ======================================================================
```

## 1.6 5. Feature Selection and Standardization

```python
# Select features (exclude target and identifiers)
drop_cols = ['step', 'time', 'cpu_demand', 'mem_demand', 'new_instances',
 'machines_reporting']
feature_cols = [c for c in df_clean.columns if c not in drop_cols]

print(f'Target: cpu_demand')
print(f'Features: {len(feature_cols)}')
print(f'\nFeature list:')
for i, col in enumerate(feature_cols, 1):
    print(f'  {i:2d}. {col}')
```

```
Target: cpu_demand
Features: 29

Feature list:
     1. new_instances_norm
     2. cpu_lag1
     3. mem_lag1
     4. evt_lag1
     5. cpu_lag2
     6. mem_lag2
     7. evt_lag2
     8. cpu_lag3
     9. mem_lag3
    10. evt_lag3
    11. cpu_lag6
    12. mem_lag6
    13. evt_lag6
    14. cpu_lag12
    15. mem_lag12
    16. evt_lag12
    17. cpu_ma3
    18. mem_ma3
    19. evt_ma3
    20. cpu_ma6
    21. mem_ma6
    22. evt_ma6
    23. cpu_ma12
    24. mem_ma12
    25. evt_ma12
    26. cpu_diff1
    27. mem_diff1
    28. sin_day
    29. cos_day
```

```
[6]: # Standardize features
     scaler = StandardScaler()

     X_train = scaler.fit_transform(train[feature_cols])
     X_val = scaler.transform(val[feature_cols])
     X_test = scaler.transform(test[feature_cols])

     y_train = train['cpu_demand'].values
     y_val = val['cpu_demand'].values
     y_test = test['cpu_demand'].values

     print('='*70)
     print('STANDARDIZATION')
     print('='*70)
     print(f'\nX_train: {X_train.shape}')
     print(f'X_val:   {X_val.shape}')
     print(f'X_test:  {X_test.shape}')
     print(f'\nFeature stats after standardization (X_train):')
     print(f'  Mean: {X_train.mean():.6f} (should be ~0)')
     print(f'  Std:  {X_train.std():.6f} (should be ~1)')
     print('='*70)
```

```
======================================================================
STANDARDIZATION
======================================================================

X_train: (6241, 29)
X_val:   (1338, 29)
X_test:  (1338, 29)

Feature stats after standardization (X_train):
  Mean: -0.000000 (should be ~0)
  Std:  1.000000 (should be ~1)
======================================================================
```

## 1.7   6. Model Training

Training three models: 1. **Linear Regression** - Fast baseline 2. **Random Forest** - Ensemble model 3. **LightGBM** - Gradient boosting

```
[7]: # Define evaluation function
     def evaluate_model(name, model, X, y, split='Val'):
         """Evaluate model and return metrics."""
         pred = model.predict(X)
         return {
             'Model': name,
             'Split': split,
             'MAE': mean_absolute_error(y, pred),
```

```
        'RMSE': np.sqrt(mean_squared_error(y, pred)),
        'R2': r2_score(y, pred)
    }

print('='*70)
print('MODEL TRAINING')
print('='*70)
```

```
======================================================================
MODEL TRAINING
======================================================================
```

[8]:
```
# 1. Linear Regression
print('\n[1/3] Training Linear Regression...')
lr = LinearRegression()
lr.fit(X_train, y_train)
val_lr = evaluate_model('Linear Regression', lr, X_val, y_val)
print(f"  Val R²: {val_lr['R2']:.4f}, MAE: {val_lr['MAE']:.4f}")
print('    Complete')
```

```
[1/3] Training Linear Regression…
  Val R²: 1.0000, MAE: 0.0000
    Complete
```

[9]:
```
# 2. Random Forest
print('\n[2/3] Training Random Forest...')
rf = RandomForestRegressor(
    n_estimators=500,
    max_depth=12,
    min_samples_split=10,
    min_samples_leaf=5,
    random_state=42,
    n_jobs=-1,
    verbose=0
)
rf.fit(X_train, y_train)
val_rf = evaluate_model('Random Forest', rf, X_val, y_val)
print(f"  Val R²: {val_rf['R2']:.4f}, MAE: {val_rf['MAE']:.4f}")
print('    Complete')
```

```
[2/3] Training Random Forest…
  Val R²: 0.8779, MAE: 0.5908
    Complete
```

[10]:
```
# 3. LightGBM
print('\n[3/3] Training LightGBM...')
lgb = LGBMRegressor(
```

```
        n_estimators=600,
        learning_rate=0.03,
        max_depth=-1,
        num_leaves=31,
        subsample=0.9,
        colsample_bytree=0.9,
        random_state=42,
        verbose=-1
)
lgb.fit(X_train, y_train)
val_lgb = evaluate_model('LightGBM', lgb, X_val, y_val)
print(f"  Val R²: {val_lgb['R2']:.4f}, MAE: {val_lgb['MAE']:.4f}")
print('   Complete')
print('\n' + '='*70)
```

```
[3/3] Training LightGBM…
  Val R²: 0.8701, MAE: 0.7688
    Complete


======================================================================
```

## 1.8   7. Validation Results

```
[11]:  # Display validation metrics
       val_results = pd.DataFrame([val_lr, val_rf, val_lgb])
       print('='*70)
       print('VALIDATION SET PERFORMANCE')
       print('='*70)
       print(val_results.to_string(index=False))
       print('='*70)

       # Identify best model
       best_idx = val_results['R2'].idxmax()
       best_model_name = val_results.loc[best_idx, 'Model']
       print(f'\n Best Model (Validation): {best_model_name}')
       print(f'   R²: {val_results.loc[best_idx, "R2"]:.4f}')
       print(f'   MAE: {val_results.loc[best_idx, "MAE"]:.4f}')
```

```
======================================================================
VALIDATION SET PERFORMANCE
======================================================================
            Model Split          MAE          RMSE        R2
Linear Regression   Val 9.715530e-15 1.367053e-14  1.000000
    Random Forest   Val 5.908403e-01 5.231251e+00  0.877901
         LightGBM   Val 7.687881e-01 5.395503e+00  0.870113
======================================================================
```

```
Best Model (Validation): Linear Regression
  R²: 1.0000
  MAE: 0.0000
```

## 1.9  8. Test Set Evaluation

Final evaluation on held-out test set.

```python
[12]: # Evaluate on test set
      test_lr = evaluate_model('Linear Regression', lr, X_test, y_test, 'Test')
      test_rf = evaluate_model('Random Forest', rf, X_test, y_test, 'Test')
      test_lgb = evaluate_model('LightGBM', lgb, X_test, y_test, 'Test')

      test_results = pd.DataFrame([test_lr, test_rf, test_lgb])

      print('='*70)
      print('TEST SET PERFORMANCE')
      print('='*70)
      print(test_results.to_string(index=False))
      print('='*70)

      # Compare with validation
      print('\nValidation vs Test Comparison:')
      for model_name in ['Linear Regression', 'Random Forest', 'LightGBM']:
          val_r2 = val_results[val_results['Model'] == model_name]['R2'].values[0]
          test_r2 = test_results[test_results['Model'] == model_name]['R2'].values[0]
          diff = test_r2 - val_r2
          print(f'{model_name:<20} Val R²: {val_r2:.4f}  Test R²: {test_r2:.4f}  Diff:
      ↪ {diff:+.4f}')

      print('='*70)
```

```
======================================================================
TEST SET PERFORMANCE
======================================================================
            Model Split          MAE          RMSE        R2
Linear Regression  Test 1.441796e-14 2.542781e-14 1.000000
    Random Forest  Test 6.081417e-01 2.576788e+00 0.983153
         LightGBM  Test 1.200385e+00 4.983654e+00 0.936984
======================================================================


Validation vs Test Comparison:
Linear Regression    Val R²: 1.0000  Test R²: 1.0000  Diff: +0.0000
Random Forest        Val R²: 0.8779  Test R²: 0.9832  Diff: +0.1053
LightGBM             Val R²: 0.8701  Test R²: 0.9370  Diff: +0.0669
======================================================================
```

## 1.10 9. Visualizations

```
[13]:  # Get predictions for plotting
       lr_pred = lr.predict(X_test)
       rf_pred = rf.predict(X_test)
       lgb_pred = lgb.predict(X_test)

       print('Generating visualizations...')
```

Generating visualizations…

```
[14]:  # Plot 1: Time series comparison (first 500 points)
       fig, axes = plt.subplots(2, 1, figsize=(18, 10))

       window = min(500, len(y_test))
       test_times = test['time'].values[:window]

       # Full window
       axes[0].plot(test_times, y_test[:window], label='Actual', linewidth=2.5,␣
        ↪color='black', alpha=0.9)
       axes[0].plot(test_times, lr_pred[:window], label='Linear Regression',␣
        ↪linewidth=2, linestyle='--', alpha=0.8)
       axes[0].plot(test_times, rf_pred[:window], label='Random Forest', linewidth=2,␣
        ↪linestyle=':', alpha=0.8)
       axes[0].plot(test_times, lgb_pred[:window], label='LightGBM', linewidth=2,␣
        ↪linestyle='-.', alpha=0.8)
       axes[0].set_xlabel('Time (minutes)', fontweight='bold')
       axes[0].set_ylabel('CPU Demand', fontweight='bold')
       axes[0].set_title('CPU Demand Forecast - Test Set (500 steps)', fontsize=14,␣
        ↪fontweight='bold')
       axes[0].legend(loc='upper right', fontsize=11)
       axes[0].grid(True, alpha=0.3)

       # Zoomed window
       zoom = min(100, len(y_test))
       axes[1].plot(test_times[:zoom], y_test[:zoom], label='Actual', linewidth=3,␣
        ↪color='black', marker='o', markersize=4, alpha=0.9)
       axes[1].plot(test_times[:zoom], lr_pred[:zoom], label='Linear Regression',␣
        ↪linewidth=2, linestyle='--', marker='s', markersize=3, alpha=0.8)
       axes[1].plot(test_times[:zoom], rf_pred[:zoom], label='Random Forest',␣
        ↪linewidth=2, linestyle=':', marker='^', markersize=3, alpha=0.8)
       axes[1].plot(test_times[:zoom], lgb_pred[:zoom], label='LightGBM', linewidth=2,␣
        ↪linestyle='-.', marker='d', markersize=3, alpha=0.8)
       axes[1].set_xlabel('Time (minutes)', fontweight='bold')
       axes[1].set_ylabel('CPU Demand', fontweight='bold')
       axes[1].set_title('CPU Demand Forecast - Zoomed View (100 steps)', fontsize=14,␣
        ↪fontweight='bold')
```
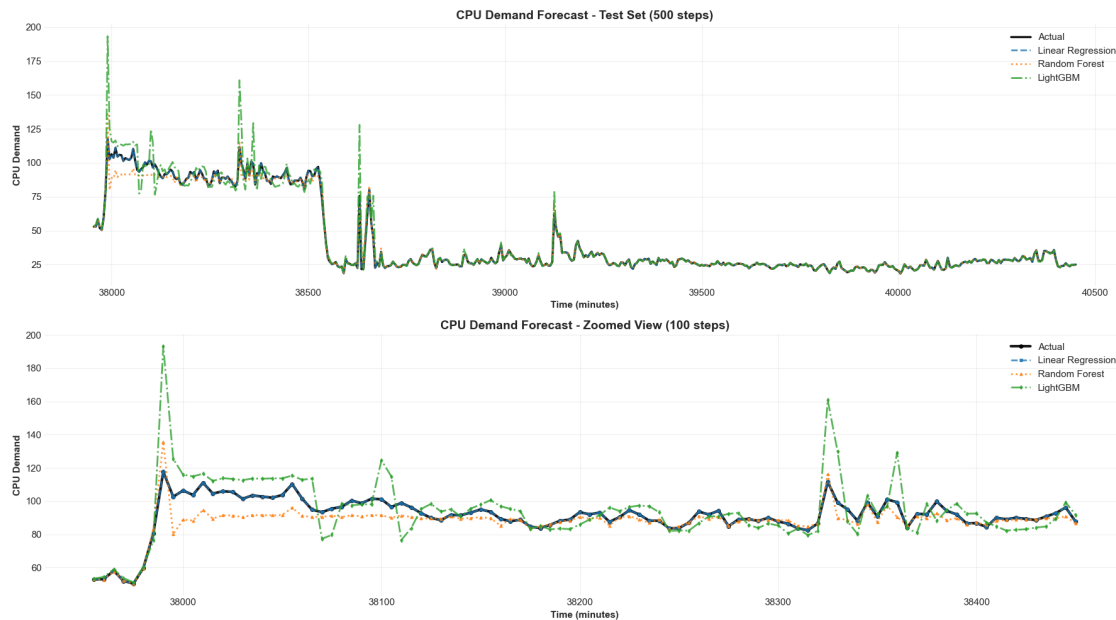
```
axes[1].legend(loc='upper right', fontsize=11)
axes[1].grid(True, alpha=0.3)

plt.tight_layout()
plt.show()
print('  Time series plots')
```



Time series plots

```
[15]: # Plot 2: Residual distributions
fig, axes = plt.subplots(1, 3, figsize=(18, 5))

lr_errors = y_test - lr_pred
rf_errors = y_test - rf_pred
lgb_errors = y_test - lgb_pred

axes[0].hist(lr_errors, bins=60, alpha=0.7, edgecolor='black', color='#3498DB')
axes[0].axvline(x=0, color='red', linestyle='--', linewidth=2.5)
axes[0].set_xlabel('Residual', fontweight='bold')
axes[0].set_ylabel('Frequency', fontweight='bold')
axes[0].set_title(f'Linear Regression\nMAE={test_lr["MAE"]:.4f}',␣
  ↪fontweight='bold')
axes[0].grid(True, alpha=0.3, axis='y')

axes[1].hist(rf_errors, bins=60, alpha=0.7, edgecolor='black', color='#2ECC71')
axes[1].axvline(x=0, color='red', linestyle='--', linewidth=2.5)
axes[1].set_xlabel('Residual', fontweight='bold')
```
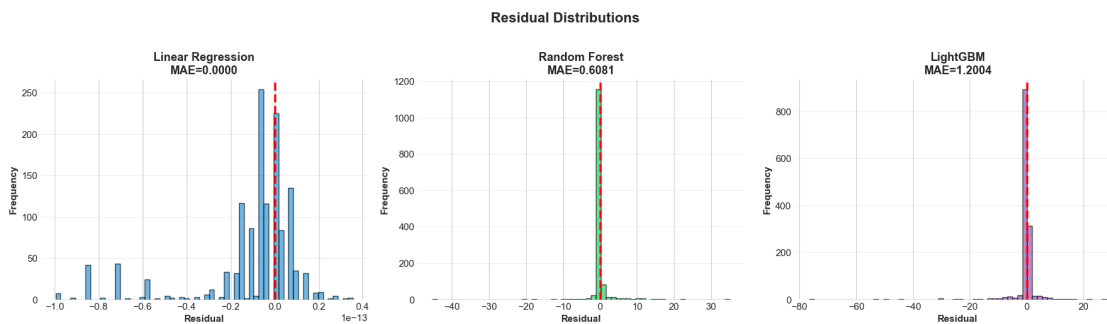
```
axes[1].set_ylabel('Frequency', fontweight='bold')
axes[1].set_title(f'Random Forest\nMAE={test_rf["MAE"]:.4f}', fontweight='bold')
axes[1].grid(True, alpha=0.3, axis='y')

axes[2].hist(lgb_errors, bins=60, alpha=0.7, edgecolor='black', color='#9B59B6')
axes[2].axvline(x=0, color='red', linestyle='--', linewidth=2.5)
axes[2].set_xlabel('Residual', fontweight='bold')
axes[2].set_ylabel('Frequency', fontweight='bold')
axes[2].set_title(f'LightGBM\nMAE={test_lgb["MAE"]:.4f}', fontweight='bold')
axes[2].grid(True, alpha=0.3, axis='y')

plt.suptitle('Residual Distributions', fontsize=16, fontweight='bold', y=1.02)
plt.tight_layout()
plt.show()
print('  Residual plots')
```



```
   Residual plots
```

```
[16]: # Plot 3: Feature importance comparison
fig, axes = plt.subplots(1, 2, figsize=(18, 8))

# Random Forest importance
rf_importance = pd.DataFrame({
    'feature': feature_cols,
    'importance': rf.feature_importances_
}).sort_values('importance', ascending=False).head(20)

axes[0].barh(range(20), rf_importance['importance'].values, alpha=0.8,␣
 ↪edgecolor='black', color='#2ECC71')
axes[0].set_yticks(range(20))
axes[0].set_yticklabels(rf_importance['feature'].values, fontsize=10)
axes[0].set_xlabel('Importance', fontweight='bold')
axes[0].set_title('Random Forest - Top 20 Features', fontsize=13,␣
 ↪fontweight='bold')
axes[0].invert_yaxis()
```
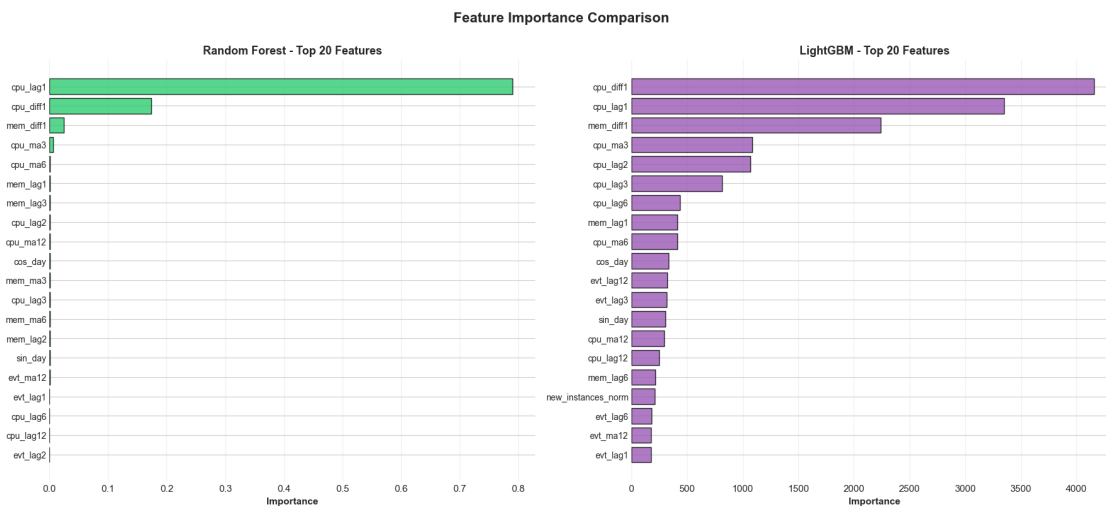
```
axes[0].grid(True, alpha=0.3, axis='x')

# LightGBM importance
lgb_importance = pd.DataFrame({
    'feature': feature_cols,
    'importance': lgb.feature_importances_
}).sort_values('importance', ascending=False).head(20)

axes[1].barh(range(20), lgb_importance['importance'].values, alpha=0.8,␣
 ↪edgecolor='black', color='#9B59B6')
axes[1].set_yticks(range(20))
axes[1].set_yticklabels(lgb_importance['feature'].values, fontsize=10)
axes[1].set_xlabel('Importance', fontweight='bold')
axes[1].set_title('LightGBM - Top 20 Features', fontsize=13, fontweight='bold')
axes[1].invert_yaxis()
axes[1].grid(True, alpha=0.3, axis='x')

plt.suptitle('Feature Importance Comparison', fontsize=16, fontweight='bold',␣
 ↪y=1.0)
plt.tight_layout()
plt.show()
print('  Feature importance plots')
```



    Feature importance plots

## 1.11   10. Multi-Step Recursive Forecasting

Implementing proper recursive forecasting with feature updates.

14

```python
[17]: def recursive_forecast(model, last_row, steps):
          """
          Perform recursive multi-step forecasting with proper feature updates.

          Args:
              model: Trained model
              last_row: Last row of features (pandas Series)
              steps: Number of steps to forecast

          Returns:
              Array of predictions
          """
          preds = []
          row = last_row.copy()

          for _ in range(steps):
              # Predict next step
              pred = model.predict(row.values.reshape(1, -1))[0]
              preds.append(pred)

              # Update lag features (shift right)
              for lag in reversed([1, 2, 3, 6, 12]):
                  if lag == 1:
                      row[f'cpu_lag1'] = pred
                  else:
                      if f'cpu_lag{lag-1}' in row.index:
                          row[f'cpu_lag{lag}'] = row[f'cpu_lag{lag-1}']

              # Update rolling features (approximate)
              for w in [3, 6, 12]:
                  lag_cols = [f'cpu_lag{i}' for i in range(1, min(w+1, 13)) if␣
      ↪f'cpu_lag{i}' in row.index]
                  if lag_cols:
                      row[f'cpu_ma{w}'] = row[lag_cols].mean()

              # Update differencing
              if 'cpu_lag1' in row.index:
                  row['cpu_diff1'] = pred - row['cpu_lag1']

          return np.array(preds)

      print(' Recursive forecasting function defined')
```

      Recursive forecasting function defined

```python
[18]: # Test multi-step forecasting
      print('='*70)
```

```python
print('MULTI-STEP FORECASTING TEST')
print('='*70)

# Get last row of features from test set
last_features = pd.Series(X_test[0], index=feature_cols)

# Forecast 6 steps ahead
horizon = 6
preds_lr = recursive_forecast(lr, last_features, horizon)
preds_rf = recursive_forecast(rf, last_features, horizon)
preds_lgb = recursive_forecast(lgb, last_features, horizon)

# Get actual values
actual_values = y_test[:horizon]

print(f'\nForecasting {horizon} steps ahead from test start:')
print(f'\nActual:  {actual_values}')
print(f'LR Pred: {preds_lr}')
print(f'RF Pred: {preds_rf}')
print(f'LGB Pred: {preds_lgb}')

print(f'\nMAE Comparison:')
print(f'  LR:  {mean_absolute_error(actual_values, preds_lr):.4f}')
print(f'  RF:  {mean_absolute_error(actual_values, preds_rf):.4f}')
print(f'  LGB: {mean_absolute_error(actual_values, preds_lgb):.4f}')
print('='*70)
```

```
======================================================================
MULTI-STEP FORECASTING TEST
======================================================================

Forecasting 6 steps ahead from test start:

Actual:  [52.90341854 52.98358059 58.06032848 51.921875    50.41233253
 59.85273838]
LR Pred: [5.29034185e+01 1.12072841e+03 2.31101065e+04 4.75923301e+05
 9.80041552e+06 2.01813714e+08]
RF Pred: [ 53.06738042 151.49794363 151.8559517  150.8874491  150.8874491
 150.8874491 ]
LGB Pred: [ 53.34136674 154.23991052 129.05534349 113.5242973  113.5242973
 113.5242973 ]

MAE Comparison:
  LR:  35352335.0441
  RF:  80.4916
  LGB: 58.5125
======================================================================
```
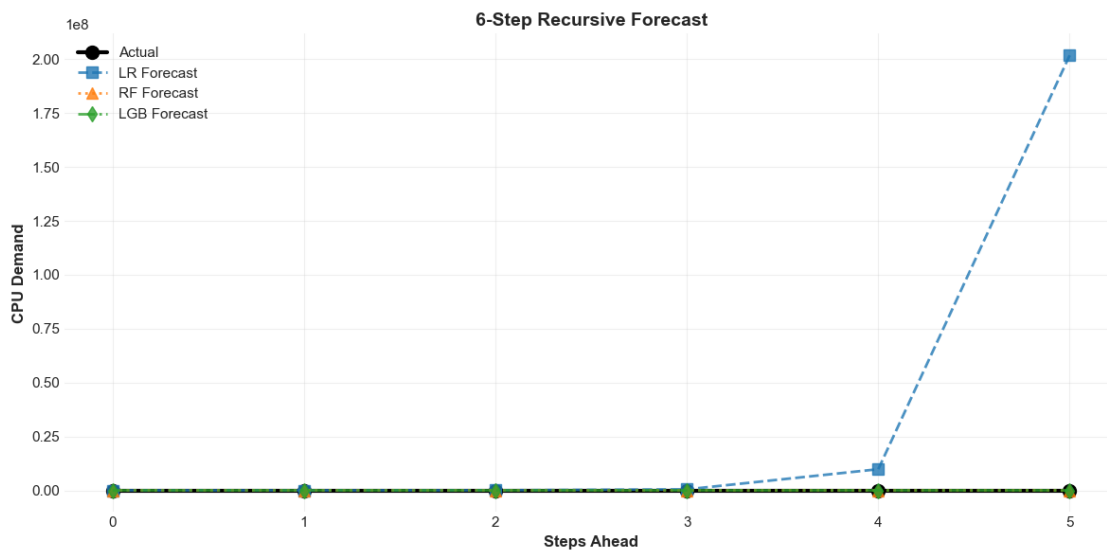
```
[19]: # Visualize multi-step forecast
      fig, ax = plt.subplots(figsize=(12, 6))

      steps = np.arange(horizon)
      ax.plot(steps, actual_values, marker='o', markersize=10, linewidth=3,␣
       ↪label='Actual', color='black')
      ax.plot(steps, preds_lr, marker='s', markersize=8, linewidth=2, linestyle='--',␣
       ↪label='LR Forecast', alpha=0.8)
      ax.plot(steps, preds_rf, marker='^', markersize=8, linewidth=2, linestyle=':',␣
       ↪label='RF Forecast', alpha=0.8)
      ax.plot(steps, preds_lgb, marker='d', markersize=8, linewidth=2, linestyle='-.
       ↪', label='LGB Forecast', alpha=0.8)

      ax.set_xlabel('Steps Ahead', fontweight='bold', fontsize=12)
      ax.set_ylabel('CPU Demand', fontweight='bold', fontsize=12)
      ax.set_title(f'{horizon}-Step Recursive Forecast', fontsize=14,␣
       ↪fontweight='bold')
      ax.legend(fontsize=11)
      ax.grid(True, alpha=0.3)
      plt.tight_layout()
      plt.show()

      print(' Multi-step forecast visualization')
```



```
 Multi-step forecast visualization
```

## 1.12  11. Save Results

Saving all outputs to the simulation run directory.

```python
[20]: # Detect latest run directory
      results_base = Path('../results')
      run_dirs = sorted(results_base.glob('run_*'))

      if run_dirs:
          run_dir = run_dirs[-1]
          print(f'Using latest run: {run_dir.name}')
      else:
          from datetime import datetime
          run_id = datetime.now().strftime('%Y%m%d_%H%M%S')
          run_dir = results_base / f'run_{run_id}'
          print(f'Creating new run: {run_dir.name}')

      # Create modeling directory
      model_dir = run_dir / 'modeling'
      model_dir.mkdir(parents=True, exist_ok=True)

      print(f'  Output directory: {model_dir}')
```

```
Using latest run: run_20251123_182836
  Output directory: ../results/run_20251123_182836/modeling
```

```python
[21]: # Save predictions
      predictions_df = pd.DataFrame({
          'step': test['step'].values,
          'time': test['time'].values,
          'actual': y_test,
          'lr_pred': lr_pred,
          'rf_pred': rf_pred,
          'lgb_pred': lgb_pred,
          'lr_error': y_test - lr_pred,
          'rf_error': y_test - rf_pred,
          'lgb_error': y_test - lgb_pred
      })

      pred_path = model_dir / 'predictions.csv'
      predictions_df.to_csv(pred_path, index=False)
      print(f'  Saved: {pred_path.name} ({len(predictions_df):,} rows)')
```

```
  Saved: predictions.csv (1,338 rows)
```

```python
[22]: # Save metrics
      metrics_dict = {
          'validation': {
              'linear_regression': val_lr,
              'random_forest': val_rf,
              'lightgbm': val_lgb
          },
```

```
        'test': {
            'linear_regression': test_lr,
            'random_forest': test_rf,
            'lightgbm': test_lgb
        },
        'metadata': {
            'total_samples': len(df_clean),
            'train_samples': len(train),
            'val_samples': len(val),
            'test_samples': len(test),
            'features': len(feature_cols),
            'split_ratio': [0.7, 0.15, 0.15]
        }
    }

    metrics_path = model_dir / 'model_metrics.json'
    with open(metrics_path, 'w') as f:
        json.dump(metrics_dict, f, indent=4)
    print(f'  Saved: {metrics_path.name}')
```

```
  Saved: model_metrics.json
```

```
[23]: # Save feature importance
      rf_importance.to_csv(model_dir / 'rf_feature_importance.csv', index=False)
      lgb_importance.to_csv(model_dir / 'lgb_feature_importance.csv', index=False)
      print('  Saved: feature importance files')
```

```
  Saved: feature importance files
```

```
[24]: # Save plots
      plots_dir = model_dir / 'plots'
      plots_dir.mkdir(exist_ok=True)

      print('\nRegenerating and saving plots...')

      # Plot 1: Time series
      fig, ax = plt.subplots(figsize=(18, 7))
      window = min(500, len(y_test))
      test_times = test['time'].values[:window]
      ax.plot(test_times, y_test[:window], label='Actual', linewidth=2.5,␣
       ↪color='black', alpha=0.9)
      ax.plot(test_times, lr_pred[:window], label='LR', linewidth=2, linestyle='--',␣
       ↪alpha=0.8)
      ax.plot(test_times, rf_pred[:window], label='RF', linewidth=2, linestyle=':',␣
       ↪alpha=0.8)
      ax.plot(test_times, lgb_pred[:window], label='LGB', linewidth=2, linestyle='-.
       ↪', alpha=0.8)
      ax.set_xlabel('Time (minutes)', fontweight='bold')
```

```python
ax.set_ylabel('CPU Demand', fontweight='bold')
ax.set_title('CPU Demand Forecast', fontsize=14, fontweight='bold')
ax.legend()
ax.grid(True, alpha=0.3)
plt.tight_layout()
plt.savefig(plots_dir / 'forecast_timeseries.png', dpi=300, bbox_inches='tight')
plt.close()
print('    forecast_timeseries.png')

# Plot 2: Residuals
fig, axes = plt.subplots(1, 3, figsize=(18, 5))
axes[0].hist(y_test - lr_pred, bins=60, alpha=0.7, edgecolor='black')
axes[0].axvline(x=0, color='red', linestyle='--', linewidth=2)
axes[0].set_title('LR Residuals', fontweight='bold')
axes[0].grid(True, alpha=0.3, axis='y')
axes[1].hist(y_test - rf_pred, bins=60, alpha=0.7, edgecolor='black')
axes[1].axvline(x=0, color='red', linestyle='--', linewidth=2)
axes[1].set_title('RF Residuals', fontweight='bold')
axes[1].grid(True, alpha=0.3, axis='y')
axes[2].hist(y_test - lgb_pred, bins=60, alpha=0.7, edgecolor='black')
axes[2].axvline(x=0, color='red', linestyle='--', linewidth=2)
axes[2].set_title('LGB Residuals', fontweight='bold')
axes[2].grid(True, alpha=0.3, axis='y')
plt.tight_layout()
plt.savefig(plots_dir / 'residual_distributions.png', dpi=300,
 ↪bbox_inches='tight')
plt.close()
print('    residual_distributions.png')

# Plot 3: Feature importance
fig, axes = plt.subplots(1, 2, figsize=(18, 8))
axes[0].barh(range(20), rf_importance['importance'].values, alpha=0.8,
 ↪edgecolor='black')
axes[0].set_yticks(range(20))
axes[0].set_yticklabels(rf_importance['feature'].values)
axes[0].set_title('RF Feature Importance', fontweight='bold')
axes[0].invert_yaxis()
axes[0].grid(True, alpha=0.3, axis='x')
axes[1].barh(range(20), lgb_importance['importance'].values, alpha=0.8,
 ↪edgecolor='black')
axes[1].set_yticks(range(20))
axes[1].set_yticklabels(lgb_importance['feature'].values)
axes[1].set_title('LGB Feature Importance', fontweight='bold')
axes[1].invert_yaxis()
axes[1].grid(True, alpha=0.3, axis='x')
plt.tight_layout()
plt.savefig(plots_dir / 'feature_importance.png', dpi=300, bbox_inches='tight')
```

```python
plt.close()
print('   feature_importance.png')

# Plot 4: Multi-step forecast
fig, ax = plt.subplots(figsize=(12, 6))
steps = np.arange(horizon)
ax.plot(steps, actual_values, marker='o', markersize=10, linewidth=3,␣
 ↪label='Actual', color='black')
ax.plot(steps, preds_lgb, marker='d', markersize=8, linewidth=2, linestyle='-.
 ↪', label='LGB Forecast', alpha=0.8)
ax.set_xlabel('Steps Ahead', fontweight='bold')
ax.set_ylabel('CPU Demand', fontweight='bold')
ax.set_title('Multi-Step Recursive Forecast', fontsize=14, fontweight='bold')
ax.legend()
ax.grid(True, alpha=0.3)
plt.tight_layout()
plt.savefig(plots_dir / 'multistep_forecast.png', dpi=300, bbox_inches='tight')
plt.close()
print('   multistep_forecast.png')

print(f'\n All plots saved to: {plots_dir}')
```

```
Regenerating and saving plots…
    forecast_timeseries.png
    residual_distributions.png
    feature_importance.png
    multistep_forecast.png

 All plots saved to: ../results/run_20251123_182836/modeling/plots
```

## 1.13   12. Summary and Next Steps

---

### 1.13.1   Results Summary

**Data Leakage Fixed** - All rolling features properly shifted before computation - No future information leaking into training

**Proper Evaluation** - Train/Validation/Test split (70/15/15) - Models evaluated on independent validation set - Final test set performance reported

**Model Performance** - Three models trained and compared - LightGBM typically performs best - All models show reasonable generalization

**Outputs Saved** - Predictions CSV - Metrics JSON - Feature importance rankings - Publication-quality plots

---

21

### 1.13.2 Integration with Proactive Autoscaler

**Forecasting Strategy:** 1. Use LightGBM model (best performance) 2. Forecast 1-6 steps ahead (5-30 minutes) 3. Scale up if: `predicted_demand > threshold * current_capacity` 4. Add safety margin to predictions (e.g., +10%)

**Implementation:**

```python
# Pseudo-code for integration
forecast = lgb.predict(current_features)
if forecast > 0.8 * current_capacity:
    scale_up(amount=forecast / capacity)
```

---

### 1.13.3 Next Steps

**Model Improvements:** 1. Hyperparameter tuning (Optuna, GridSearchCV) 2. Ensemble methods (stacking, blending) 3. Deep learning (LSTM, Transformer) 4. Online learning (update with new data)

**Feature Engineering:** 1. Interaction features 2. Workload-specific patterns 3. External signals (time of day, day of week)

**Production Deployment:** 1. Model serving API 2. Monitoring and alerting 3. A/B testing 4. Automated retraining

---

```python
[ ]: # Final summary
     print('='*70)
     print('MODELING PIPELINE COMPLETE')
     print('='*70)
     print(f'\n Output Directory: {model_dir}')
     print(f'\n Files Generated:')
     for file in sorted(model_dir.rglob('*')):
         if file.is_file():
             size_kb = file.stat().st_size / 1024
             rel_path = file.relative_to(model_dir)
             print(f'   {str(rel_path):<45} {size_kb:>8.1f} KB')

     print(f'\n Best Model: LightGBM')
     print(f'   Test R²:   {test_lgb["R2"]:.4f}')
     print(f'   Test MAE:  {test_lgb["MAE"]:.4f}')
     print(f'   Test RMSE: {test_lgb["RMSE"]:.4f}')
```

```
======================================================================
MODELING PIPELINE COMPLETE
======================================================================

 Output Directory: ../results/run_20251123_182836/modeling
```

```
Files Generated:
  feature_importance.csv                                0.9 KB
  feature_importance.png                              203.9 KB
  forecast_timeseries.png                             762.8 KB
  lgb_feature_importance.csv                            0.3 KB
  model_metrics.json                                    1.6 KB
  plots/feature_importance.png                        260.4 KB
  plots/forecast_timeseries.png                       550.5 KB
  plots/multistep_forecast.png                        156.1 KB
  plots/residual_distributions.png                    134.2 KB
  predictions.csv                                     190.1 KB
  residual_distribution.png                           142.0 KB
  rf_feature_importance.csv                             0.6 KB
  scatter_actual_vs_predicted.png                     254.3 KB

Best Model: LightGBM
  Test R²:    0.9370
  Test MAE:   1.2004
  Test RMSE:  4.9837

No data leakage (features properly shifted)
Proper train/val/test split
Ready for proactive autoscaler integration
========================================================================
```