# JAX Machine Learning Framework Presentation

## Executive Summary

This presentation by Siwat (Customer Solutions Consultant at Google Thailand) provides a comprehensive introduction to JAX, Google's high-performance Python library for machine learning. The talk covers JAX's core concepts, performance advantages, and practical applications in modern AI development.
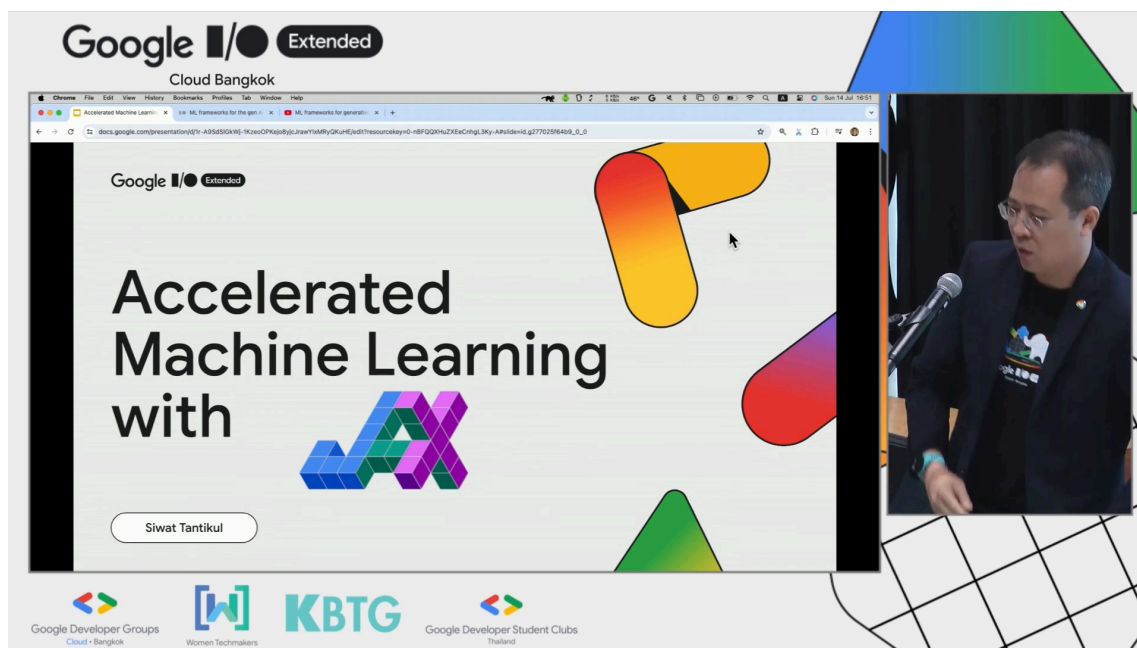
## Speaker Highlights

**Siwat** - Google Customer Solutions Consultant

- Extensive experience with JAX implementation at Google
- Covers JAX usage from AlphaFold to Gemini
- Provides practical insights and real-world performance comparisons

## Timeline Summary

### Introduction & Context (0-300s)



The presentation begins with an introduction to JAX and its significance in modern AI development. JAX has been used extensively at Google, from AlphaFold protein prediction to current Gemini models. The speaker emphasizes the growing computational requirements - from 10^15 petaflops for Transformer models in 2017 to over 20,000 billion petaflops for models like GPT-4.

### Pure Functions & Functional Programming (375-465s)

A critical concept in JAX is pure functional programming. Unlike traditional Python functions that may mutate input data, JAX requires immutable inputs. The speaker demonstrates how pure functions create copies rather than modifying original data, enabling better caching and optimization.

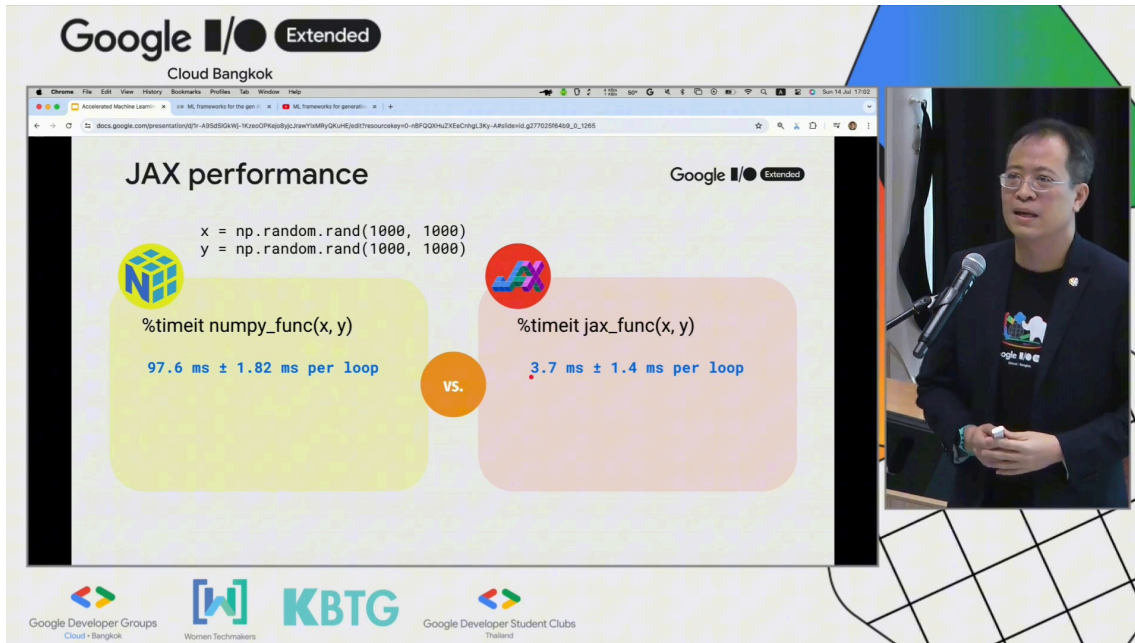## JAX Stack Architecture (462-646s)



The JAX stack consists of:

- Hardware layer (CPU/GPU/TPU)
- XLA (Accelerated Linear Algebra) compiler
- JAX core library

- Ecosystem libraries (Flax, Optax, etc.)

XLA is the heart of JAX, converting Python code into optimized machine code for hardware acceleration.

## Performance Comparison (649-780s)



Dramatic performance improvements demonstrated:

- NumPy: 97 milliseconds
- JAX: 3 milliseconds
- JAX with JIT: 146 microseconds

The speaker emphasizes that JAX excels with large-scale vectorized operations.

## JIT Compilation (695-860s)

Just-In-Time (JIT) compilation is explained as a key performance booster. The first execution compiles the function, subsequent calls use cached compiled code. Simple decorator syntax ( `@jit` ) makes it easy to apply.

### Auto-Differentiation (895-1019s)



JAX's automatic differentiation capabilities are showcased. The `jax.grad` function automatically computes derivatives regardless of equation complexity, supporting nested gradients of any degree. Practical example shown with logistic regression.

### Vectorization with vmap (1019-1175s)

Vector mapping (vmap) enables parallel operations across data batches:

- Traditional for-loop: 7 seconds
- With vmap: 1 second
- With vmap + JIT: 132 microseconds

This demonstrates massive speedups for batch processing operations.

### XLA Internals (1180-1302s)



Deep dive into how JAX converts Python code through multiple compilation stages:

1. JAX expressions (jaxpr)
2. StableHLO intermediate representation
3. XLA optimization
4. Machine code for hardware

## Neural Networks & MaxText (1411-1615s)



JAX ecosystem for neural networks includes Flax and Equinox libraries. MaxText is introduced as Google's open-source framework for training large language models, successfully used to train models with 50,000 TPUs achieving 55% model FLOP utilization.

## Getting Started Resources (1640-1776s)

Recommended learning resources:

- JAX tutorials on Kaggle (9 lessons)
- Awesome JAX repository
- JAX ReadTheDocs official documentation
- Google's JAX GitHub
- Book by Grigory Sapunov for deep understanding

## Key Takeaways

- JAX provides NumPy-compatible syntax with dramatic performance improvements
- Core features: JIT compilation, auto-differentiation, vectorization, parallelization
- Used extensively at Google for production AI systems
- Open-source with strong ecosystem (Flax, Optax, MaxText)
- Ideal for high-performance machine learning and large-scale model training
- Functional programming paradigm enables better optimization and caching
- Performance gains of 100-1000x possible with proper usage of JIT and vmap