# Generative AI: AI Intelligence at Cloud Scale

**Unlock the Future with Generative AI**

Ben Rodrigue - Generative AI Product Leader - BenRod@TensorIoT.com
Yuxin Yang - Machine Learning Practice Manager - Yuxin.Yang@TensorIoT.com

1

# Making Things Intelligent.

## World Class IoT, Data, AI & ML on AWS

■ Cloud Native Products and Consulting

■ Headquartered in Southern California, Global Presence

■ Focus on IoT, AI, ML, Data, App Development & Modernization

■ Meet the customer wherever they are in their cloud adoption journey

■ Named Sustainability Partner of the year at re:Invent 2022

aws

**PARTNER**
Advanced Tier
Services

▪ ML Services Competency
▪ IoT Services Competency
▪ Retail Services Competency
▪ Industrial Software Service
  Competency
▪ Travel &
  Hospitality Services

2022
**AWS Sustainability Partner of the Year**
North America

TensorIoT
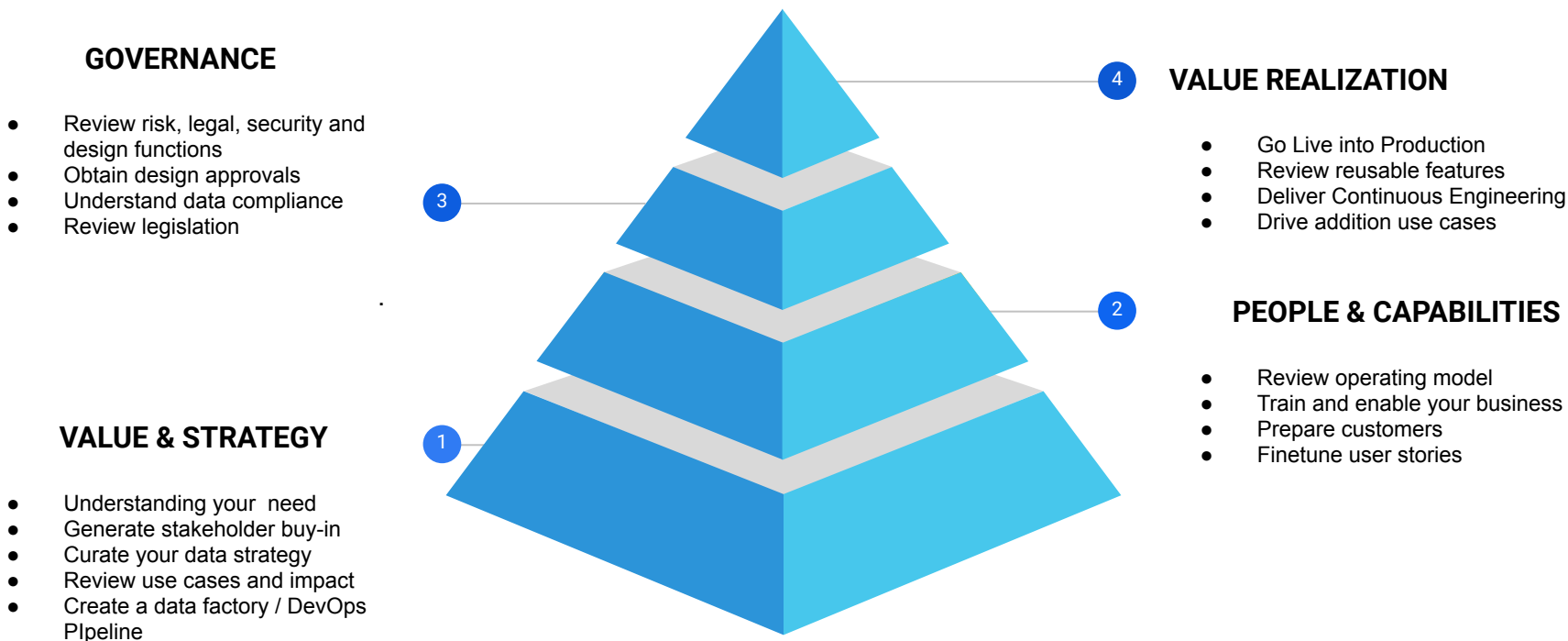
2

TensorIoT

# What is Generative AI?

**TensorIoT**

What if you could feed a Generative AI Model with your company's data to automate business functions, generate content, or drive critical decisions all while retaining control of your data?

# Generative AI:
# Adoption Journey with TensorIoT

**Unlock the Future with Generative AI**

# Generative AI adoption journey

## We meet you where you are.

**6**

**GOVERNANCE**

- Review risk, legal, security and design functions
- Obtain design approvals
- Understand data compliance
- Review legislation

**3**

**VALUE REALIZATION**

**4**

- Go Live into Production
- Review reusable features
- Deliver Continuous Engineering
- Drive addition use cases

**PEOPLE & CAPABILITIES**

**2**

- Review operating model
- Train and enable your business
- Prepare customers
- Finetune user stories

**VALUE & STRATEGY**

**1**

- Understanding your need
- Generate stakeholder buy-in
- Curate your data strategy
- Review use cases and impact
- Create a data factory / DevOps PIpeline

TensorIoT

# Results at Velocity

**We provide tangible results in weeks, not years**

| Executive Briefing | AI Assessment | Proof of Value | Deploy Model into Production |
|---|---|---|---|
| 2 hours | 5 days | 2-4 weeks | 2-3 months |

TensorIoT

# Results at Velocity

**We provide tangible results in weeks, not years**

**Executive Briefing: (2 hrs)**
Get insights in TensorIoT's AI best practices, capabilities, and high value use cases

**Generative AI Assessment: (5 days)**
Assess your organization's readiness to embark on a Generative AI journey and build a roadmap to deploy internal AI services.

**Proof of Value: (2-4 weeks)**
Identify high impact use case and rapidly test a foundational model against your data set.

**Deploy Models into Production:(2-3 months)**
Train, Refine, Deploy, and Scale your model into production and begin reaping the benefits and economic value of Generative AI.

TensorIoT

TensorIoT

# Generative AI: Strategy Assessment

**Unlock the Future with Generative AI**

# Generative AI: Readiness Assessment

## Objectives and Benefits

**Objective:**

Assess your organization's readiness to embark on a Generative AI journey and build a roadmap to deploy internal AI services.

**Benefits:**

- Prepare for the business disruption from Foundation Models, which are seen in action with technologies like ChatGPT, BARD, Amazon Bedrock & Hugging Face.
- Determine Generative AI's potential value in your business context.
- Assess your organization's readiness for AI implementation.
- Identify AI use cases aligned with your strategic goals.
- Ensure secure, unbiased, and compliant AI deployment.
- Develop a customized roadmap for seamless AI adoption.

TensorIoT

# Assessment Overview:

- **Duration:** 5 days with 2-4 hour sessions
- **Participants:** Cross-functional teams from IT, Business, Data Science, and Security
- **Output:** A tailored AI readiness assessment scorecard, use case catalogue, and roadmap to guide your Generative AI journey

| Introduction/Discovery | Discuss Strategies | Develop Roadmap | Outcome |
| --- | --- | --- | --- |
| **Business Overview:** | **Strategy and Capabilities:** | **Recommend Next Steps:** | **Outcome:** |
| Introduction to Generative AI, assessing your organization's resources and capabilities, identifying high-priority goals, and exploring potential Generative AI use cases. | Reviewing your current cross-organizational strategy to leverage Generative AI, reviewing current data footprint maturity on AWS, and exploring the long-term strategy for Generative AI in your organization. | Developing a tailored roadmap for Generative AI maturity and adoption, identifying potential challenges and mitigation strategies, and discussing the next steps to achieve your Generative AI objectives. | Determine an action plan for roadmap implementation, encompassing design choices, AWS architecture, timeline projections, and the identification of high-impact use cases to achieve both short-term and long-term objectives. |

TensorIoT

# Generative AI: Proof of Value

**Unlock the Future with Generative AI**

# Generative AI: Proof of Value

## Objectives and Benefits

### Overview:

**Wizdom AI** is an all-inclusive technical framework, enabling you to train and deploy Generative AI, AWS Services, and Foundational Models like Bloom, Titan, and Jurassic within your own AWS account.

### Objective:

Discover a high-impact use case and swiftly evaluate a foundational model using your data set, showcasing immediate value to all stakeholders.

### Benefits:

- Functional PoV in 2-4 weeks
- Stays in your account giving you full control to the model and your data
- Deployed using out of the box, mature, Foundational Models
- Build on AWS leveraging next-gen services

**Wizdom AI**

TensorIoT

# Generative AI: Proof of Value

## Solution Highlights

| | Value | 3rd party (OpenAI) | Wizdom AI Framework | Additional Details |
|---|---|---|---|---|
| 1 | You control the Model | ✕ | ✓ | • All Models are deployed in your AWS account<br>• Foundational Models are Open Source or licensed by you |
| 2 | You control your Data | ✕ | ✓ | • Full Privacy<br>• Your data never leaves your account |
| 3 | End to End Security | ✕ | ✓ | • Maintain industry compliance<br>• Maintain data sovereignty<br>• Follows security best practices |
| 4 | Fully integrated with AWS services | ✕ | ✓ | • Seamlessly integrates with AWS services<br>• On demand, scalable managed services<br>• Cost optimized |
| 5 | Dedicated Support and Maintenance | ✕ | ✓ | • Support and maintenance available<br>• Ongoing model performance available<br>• Ongoing updates available |

TensorIoT
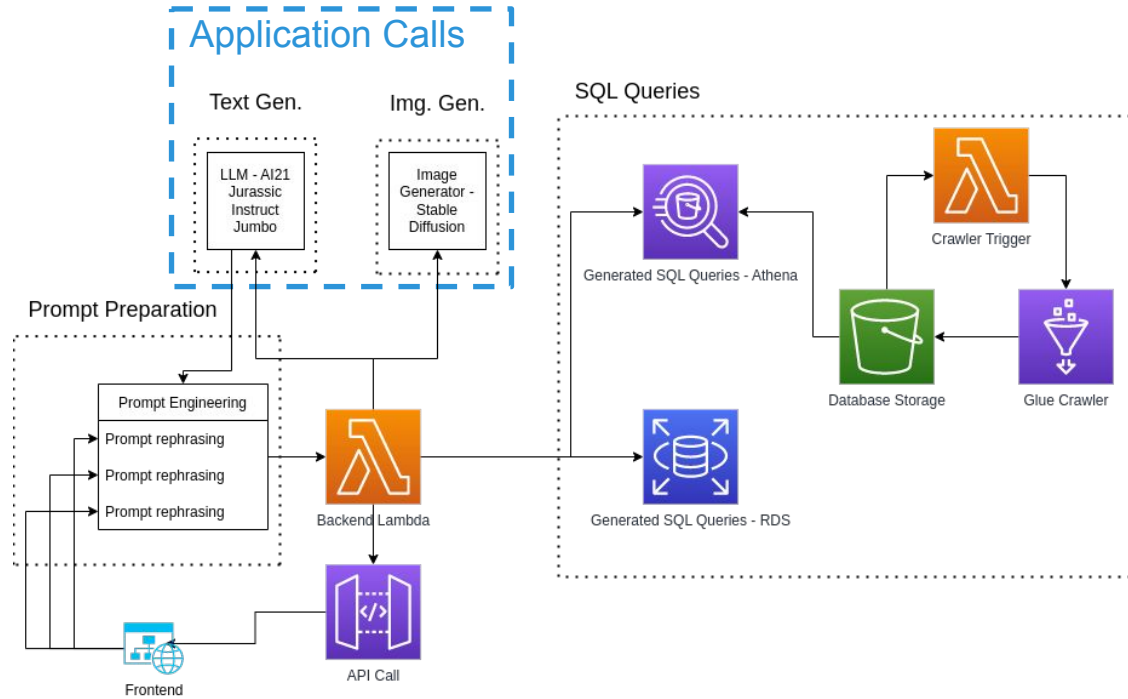
# Generative AI: Proof of Value

**Execution Journey**

Design

Solution Build

Week 1

Week 2

Week 3

Week 4

Solution Build

Test
Demo
Handoff

TensorIoT

TensorIoT

# Example Architectures

# Wizdom AI: Inference Backend

**Customer's own VPC environment**

Model artifacts stored in an Amazon S3 bucket

Amazon ECR container image

Create an Amazon SageMaker Model

Choose a SageMaker Inference option

(Optional: Use Inference Recommender to get instance recommendations)

Create an Amazon SageMaker Endpoint Configuration

Create an Amazon SageMaker Endpoint

**Application Calls**

Invoke endpoint to get inferences

AI21 labs

amazon

# Wizdom AI: Core Architecture

**Customer's own VPC environment**

### Application Calls

**Text Gen.**

LLM - AI21 Jurassic Instruct Jumbo

**Img. Gen.**

Image Generator - Stable Diffusion

**SQL Queries**

Generated SQL Queries - Athena

Crawler Trigger

Database Storage

Glue Crawler

**Prompt Preparation**

Prompt Engineering
Prompt rephrasing
Prompt rephrasing
Prompt rephrasing

Backend Lambda

Generated SQL Queries - RDS

Frontend

API Call

# Wizdom AI: Fine-tuning LLM Architecture

Customer's own VPC environment

Data Sources

Customizable Training Data

LLM Model containers

SageMaker Studio Environment

Sagemaker Notebook

Sagemaker Training Jobs

Model artifacts

Sagemaker Model Registry

# Transformers Deployed On AWS

**Hugging Face**
10,000+ pre-trained Hugging Face Transformers models for NLP, speech, and vision models

## Build

Develop scripts on Amazon SageMaker Notebook Instances, SageMaker Studio, or on your IDE

## Train

Train in Hugging Face Deep Learning Containers (DLC)

Fine-tune and manage experiments

## Deploy

Deploy any Hugging Face model easily on Amazon SageMaker

Automatically monitor and scale model endpoints

Download model from Amazon S3 for self-managed deployment

TensorIoT

# Thank you

Ben Rodrigue - Generative AI Product Leader - BenRod@TensorIoT.com
Yuxin Yang - Machine Learning Practice Manager - Yuxin.Yang@TensorIoT.com

# Wizdom AI: Model Training

**ML Frameworks**       PyTorch      TensorFlow

**AWS SageMaker Training Platform**

Large-scale Clusters    Fault-tolerant    Hugging Face Integration    SageMaker Training Compiler

**AWS SageMaker Training Parallelism**

**Model Parallelism**      **Data Parallelism**

TensorIoT

# AWS FM deployment in Production (Deployment + Inference)

```
In [28]:    region = boto3.Session().region_name
            if region not in model_package_map.keys():
                raise ("UNSUPPORTED REGION")

            model_package_arn = model_package_map[region]
```

```
In [29]:    role = get_execution_role()
            sagemaker_session = sage.Session()

            runtime_sm_client = boto3.client("runtime.sagemaker")
```

## 2. Create an endpoint and perform real-time inference

If you want to understand how real-time inference with Amazon SageMaker works, see Documentation.

```
In [30]:    endpoint_name = "j1-grande"

            content_type = "application/json"

            real_time_inference_instance_type = (
                "ml.g5.12xlarge"
            )
```

### A. Create an endpoint

```
In [31]:    # create a deployable model from the model package.
            model = ModelPackage(
                role=role, model_package_arn=model_package_arn, sagemaker_session=sagemaker_session
            )

            # Deploy the model
            predictor = model.deploy(1, real_time_inference_instance_type, endpoint_name=endpoint_name,
                                     model_data_download_timeout=3600,
                                     container_startup_health_check_timeout=600,
                                     )
            ---------------!
```

## Foundation models

| Search for a model |

### AI21 Summarize
By AI21 Labs | Ver 1.1.000

**THE INPUT TEXT SHOULD CONTAIN AT LEAST 40 *WORDS* AND NO MORE THAN 50,000 *CHARACTERS*. THIS TRANSLATES TO ROUGHLY 10,000 WORDS, OR AN IMPRESSIVE 40 PAGES!**

Summarize texts with our world-class summarization engine. Quick integration with high quality for all kinds of text.

[ View model ]

### AI21 Jurassic-2 Jumbo
By AI21 Labs | Ver 1.0.032

**PRE-TRAINED LANGUAGE MODEL TRAINED BY AI21 LABS ON A CORPUS OF WEB TEXT INCLUDING NATURAL LANGUAGE AND COMPUTER PROGRAMS WITH RECENT DATA - UPDATED TO MID 2022. THIS MODEL HAS A 8192 TOKEN CONTEXT WINDOW (I.E. THE LENGTH OF THE PROMPT + COMPLETION SHOULD BE AT MOST 8192 TOKENS).**

Best-in-class large language model designed for maximum quality. Ideal for generating text using plain English.

[ View model ]

### AI21 Jurassic-2 Jumbo Instruct
By AI21 Labs | Ver 1.1.033

**PRE-TRAINED LANGUAGE MODEL TRAINED BY AI21 LABS ON A CORPUS OF WEB TEXT INCLUDING NATURAL LANGUAGE AND COMPUTER PROGRAMS WITH RECENT DATA - UPDATED TO MID 2022. THIS MODEL HAS A 8192 TOKEN CONTEXT WINDOW (I.E. THE LENGTH OF THE PROMPT + COMPLETION SHOULD BE AT MOST 8192 TOKENS).**

Best-in-class instruction following model designed for maximum quality. Ideal for generating text using plain instructions.

[ View model ]

TensorIoT

# Hugging Face in Production (Deployment + Inference)

```python
role = sagemaker.get_execution_role()
# Hub Model configuration. https://huggingface.co/models
hub = {
    'HF_MODEL_ID':'tscholak/cxmefzzi',
    'HF_TASK':'text2text-generation'
}

# create Hugging Face Model Class
huggingface_model = HuggingFaceModel(
    transformers_version='4.17.0',
    pytorch_version='1.10.2',
    py_version='py38',
    env=hub,
    role=role,
)

# deploy model to SageMaker Inference
predictor = huggingface_model.deploy(
    initial_instance_count=1, # number of
    instance_type='ml.g4dn.4xlarge' # ec2
)

---------------!
```

```
In [8]: predictor.predict({
            'inputs': "How many singers do we have? \
            | concert_singer \
            | stadium : stadium_id, location, name, capacity, highest, lowest, average \
            | singer : singer_id, name, country, song_name, song_release_year, age, is_male \
            | concert : concert_id, concert_name, theme, stadium_id, year | singer_in_concert : concert_id, singer_id"
        })

Out[8]: [{'generated_text': 'concert_singer | select count(*) from singer'}]
```

```
In [9]: predictor.predict({
            'inputs': "What are the unique singers that have had any concerts? \
            | concert_singer | stadium : stadium_id, location, name, capacity, highest, lowest, average \
            | singer : singer_id, name, country, song_name, song_release_year, age, is_male \
            | concert : concert_id, concert_name, theme, stadium_id, year | singer_in_concert : concert_id, singer_id"
        })

Out[9]: [{'generated_text': 'concert_singer | select count(distinct singer_id) from singer_in_concert'}]
```

```
In [10]: predictor.predict({
             'inputs': "Select all the unique singers that have had any concerts? \
             | concert_singer \
             | stadium : stadium_id, location, name, capacity, highest, lowest, average \
             | singer : singer_id, name, country, song_name, song_release_year, age, is_male \
             | concert : concert_id, concert_name, theme, stadium_id, year \
             | singer_in_concert : concert_id, singer_id"
         })

Out[10]: [{'generated_text': 'concert_singer | select distinct singer_id from singer_in_concert'}]
```

Search Relevant Information

Knowledge Sources

2 Query

Relevant Information for Enhanced Context

3

1 Prompt + Query

Generated Text Response 5

Prompt + Query + Enhanced Context
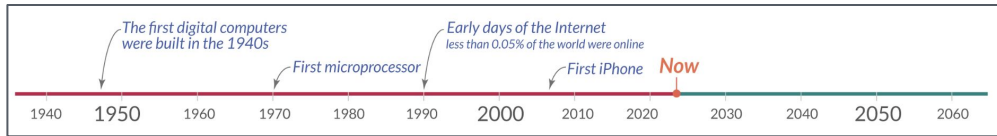
4

Large Language Model EndPoint

# Trends in Artificial Intelligence

**AI Index Annual Report 2023**

1. **Industry races ahead of academia**
2. **Performance saturation on traditional benchmarks.**
3. **AI is both helping and harming the environment.**
4. **The world's best new scientist … AI?**
5. **The number of incidents concerning the misuse of AI is rapidly rising.**
6. **The demand for AI-related professional skills is increasing across virtually every American industrial sector.**
7. **While the proportion of companies adopting AI has plateaued, the companies that have adopted AI continue to pull ahead.**
8. **Policymaker interest in AI is on the rise.**
9. **Chinese citizens are among those who feel the most positively about AI products and services. Americans … not so much.**

SOURCE: https://aiindex.stanford.edu/report/

# History of AI





https://ourworldindata.org/brief-history-of-ai

# Scale - Good / Bad

Here are some of the potential benefits of generative AI:

- It can be used to create new and innovative content.
- It can be used to personalize content for individual users.
- It can be used to automate tasks that are currently done by humans.
- It can be used to generate new insights and ideas.

Here are some of the potential risks of generative AI:

- It could be used to create fake news and propaganda.
- It could be used to generate harmful or offensive content.
- It could be used to automate jobs that are currently done by humans.
- It could lead to job losses and economic disruption.

# Generative AI: Advisory Services

**Unlock the Future with Generative AI**

# Generative AI: Advisory Services
## Objectives and Benefits

### Overview:

**Wizdom AI** is an all-inclusive technical framework, enabling you to train and deploy Generative AI, AWS Services, and Foundational Models like Bloom, Titan, and Jurassic within your own AWS account.

### Objective:

Discover a high-impact use case and swiftly evaluate a foundational model using your data set, showcasing immediate value to all stakeholders.

### Benefits:

- Functional PoV in 2-4 weeks
- Stays in your account giving you full control to the model and your data
- Deployed using out of the box, mature, Foundational Models
- Build on AWS leveraging next-gen services

**Wizdom AI**