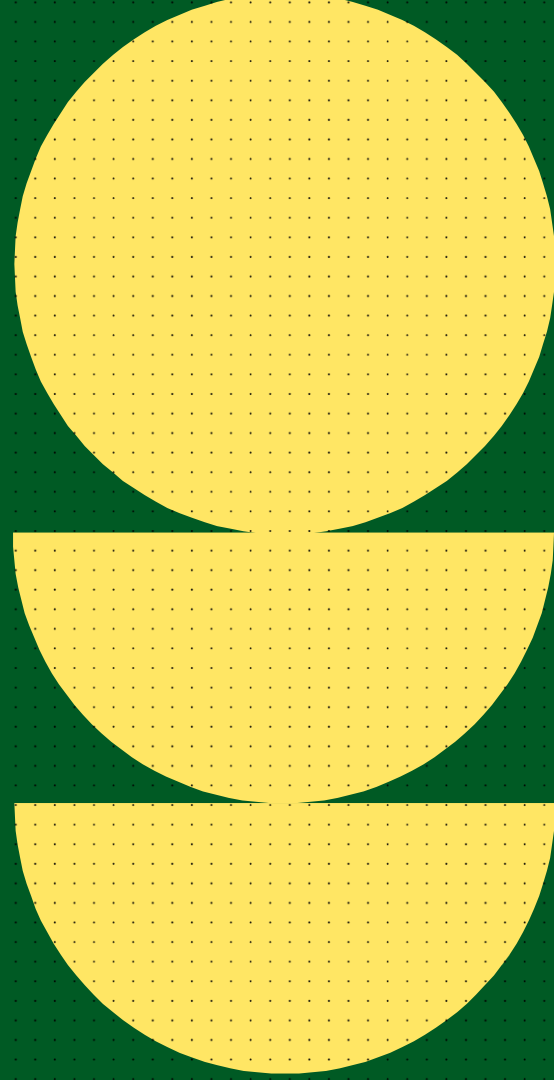# Expedition To The LLM & AI Frontier

**Sanghamitra Goswami**
Senior Director, PagerDuty
Data Science & Machine Learning
**October 17th, LeadDev West Coast 2023**

**Introduction to PagerDuty**

- Pagerduty Operations Cloud
- PagerDuty - AI
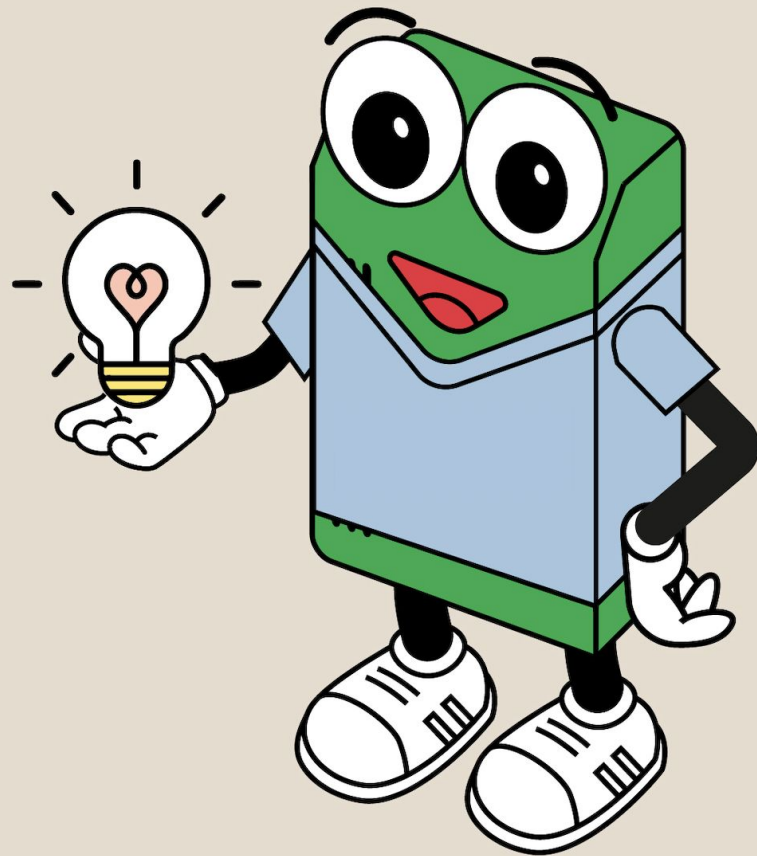
**AI to Generative AI**

- Challenges
- Multi-modal Approach
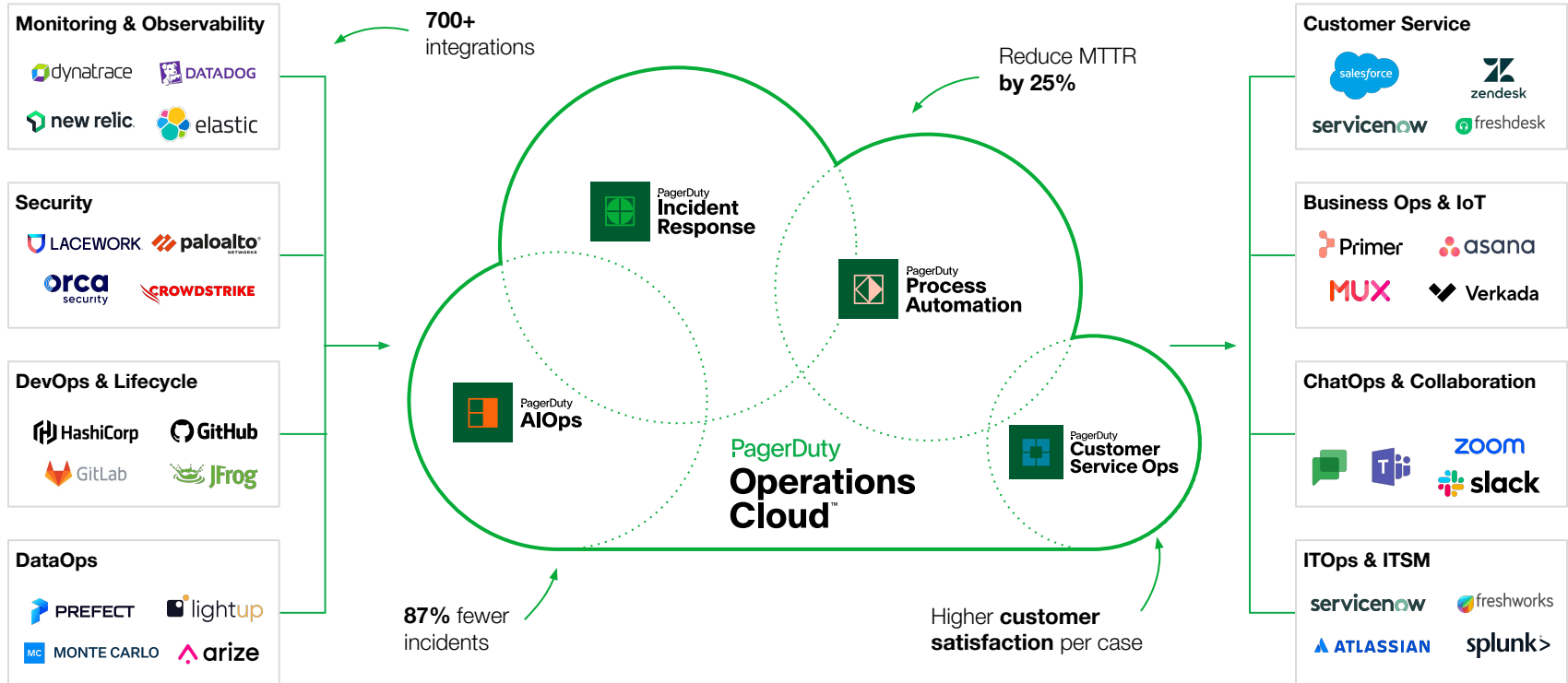
**New Generative AI Features**

- SMART STUNT
- AI Generated Postmortems
- AI Generated Runbooks

# PagerDuty

Through its SaaS-based platform, PagerDuty empowers developers, DevOps, IT operations and business leaders to prevent and resolve business-impacting incidents

# PagerDuty Operations Cloud Ecosystem

**Monitoring & Observability**
dynatrace
DATADOG
new relic
elastic

**Security**
LACEWORK
paloalto NETWORKS
ORCA security
CROWDSTRIKE

**DevOps & Lifecycle**
HashiCorp
GitHub
GitLab
JFrog

**DataOps**
PREFECT
lightup
MONTE CARLO
arize

**700+** integrations

**PagerDuty Incident Response**

**PagerDuty Process Automation**

**PagerDuty AIOps**

**PagerDuty Operations Cloud™**

**PagerDuty Customer Service Ops**

Reduce MTTR **by 25%**

**87%** fewer incidents

Higher **customer satisfaction** per case

**Customer Service**
salesforce
zendesk
servicenow
freshdesk

**Business Ops & IoT**
Primer
asana
MUX
Verkada

**ChatOps & Collaboration**
zoom
slack

**ITOps & ITSM**
servicenow
freshworks
ATLASSIAN
splunk>

PagerDuty

# AI is part of our DNA

AI has been part of our platform



**AIOps**

**Incident Response**

**Process Automation**

**Customer Service Ops**

**We plan to incorporate AI throughout our product portfolio**
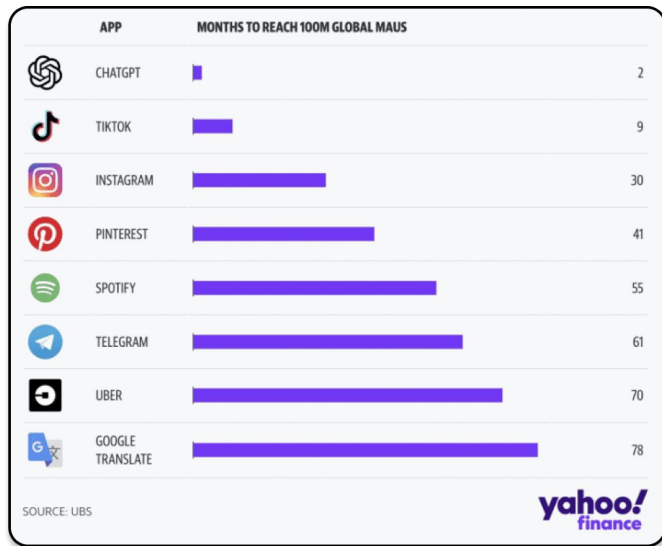
PagerDuty

# PagerDuty AI

Reduce Noise, Create Context, Automate Toil

**Suppression & Deduplication**

**Automation Actions**

**INGEST**

| DETECT | PREVENT | MOBILIZE | DIAGNOSE | RESOLVE |
|---|---|---|---|---|

**Alert Grouping**
**Auto-Pause Incidents**

**Related Incidents**
**Past Incidents**
**Probable Origin**
**Outlier Incidents**
**Change Correlation**

PagerDuty

# AI To Generative AI - Our Expedition

# Generative AI Landscape a Few Months Back



| APP | MONTHS TO REACH 100M GLOBAL MAUS | |
|-----|----------------------------------|---|
| CHATGPT | | 2 |
| TIKTOK | | 9 |
| INSTAGRAM | | 30 |
| PINTEREST | | 41 |
| SPOTIFY | | 55 |
| TELEGRAM | | 61 |
| UBER | | 70 |
| GOOGLE TRANSLATE | | 78 |

SOURCE: UBS

yahoo! finance

**10 months ago** OpenAI released ChatGPT
- Crazy adoption rate as people were amazed by its capabilities.
- Sparked a lot of businesses to integrate it into their products
- We wanted to innovate and add to our Product

Our Goal:

Let's adopt Generative AI and develop Product features

PagerDuty

# Key Challenges

- New planning
- Ratio of data scientists to software engineers low
- New technology
- Concurrent workflows
- Roadmap realignment

PagerDuty

# Our Multi-Modal Approach

- Empower & Educate
- Establish Rules of Engagement
- Discover Your Use Cases
- Establish Data Privacy & Legal guidelines
- Determine the LLM of Your Choice
- Set Operating Principles To Release

# Empower & Educate

Hack Week → Work in a team, learn from each other

Several Gen AI focussed slack channels → Share the learning, curated courses

Data Science/AI Office hours → Ask questions to your Data Science Team

Form a core lead group→ A multidisciplinary group with AI, legal, security and engineering

Dedicate x% of resources from all engineering teams to learn → AI champions in every team

# Establish Rules on Engagement

- Involve a Data Scientist in early conversations
- Start a design document early
- Involve Legal and Security early



PagerDuty

# You Need Data Scientists

# Role Evolving…

Explore data feasibility                                    Garbage in garbage out

Building prompt Vs. Fine tuning                             Costly interactions

Build within a Generative AI feature                        Make Product unique

Evaluate feature performance                                Measure performance & quality

Assist Product in selecting right customers                 Product customer fit

PagerDuty

# Discover your Use Cases

**Code Generation**



- We pivoted to one of our values **#champinionthecustomer**
- Where can Gen AI help?
- What are some of the use cases for PagerDuty?
- Limitations?

**Chat Bot**



**Summarization**



PagerDuty

# Determine the LLM of your Choice

Accuracy  → Score outputs from LLMs

Cost → Queries, especially those requiring long prompts can be expensive

Latency → It is not uncommon for a query to take several seconds

Context Window → Several of our use cases require large context windows

Data Governance → Control over our data

Ease of setup & SLA → Effort for PagerDuty to access, manage and maintain

# Operating Principles To Release

- Internal Preview

- External Preview

- Early Access

- General Access

- Human-in-the-loop

# Key Takeaways

- Preserve Optionality
  - Build for evolution, not perfection
- Manage Risks
  - Customer Trust is the #1 priority
- Focus on customer adoption
  - Show new value against their jobs-to-be-done

# GenAI features just announced

**1** AI-Generated
Status Updates

**2** AI-Generated
Incident Postmortems

**3** AI-Generated
Process Automation

# One-click status updates, built for any audience

– Generate persona-specific status updates with just a few clicks

– Make it easy for responders to keep internal stakeholders and executives in the loop

Demo

**Waitlist available**

# Save time and remove toil from postmortems

- Use AI to create a comprehensive summary of what happened, when, how it was resolved, and key actions for next time

- Automatically collect and collate incident data (including logs, metrics, and relevant Slack or Microsoft Teams conversations

Demo

**Waitlist available**

# Generate Process Automation job definitions via AI

– Get a fast start in developing new workflows

– Speed up learning and problem solving

– Generate results for any scripting language

Demo

**Waitlist available**

# What's Next

- Sign up for PagerDuty's waitlist
- Watch an extended demo that shows some of these features in action https://youtu.be/OgBS2cvkiCs
- Hear more from the team on learnings https://youtu.be/2mEVY6rmX3M

**PagerDuty**

Thank You

# Establish Rules on Engagement

**Start a design document early**

- This tightens everything up and gives a basis for discussion
- Include the API contract when working with external vendors

**Involve Legal and Security Early**

- We partnered with legal from the beginning, involving them in the earliest stages of the process. This allowed us to get early advice on how to think about the short, medium and long-term release phases

Presented in DataCon LA 2023 by Sanghamitra Goswami

# Key Challenges

**New planning** → We wanted to move fast while the industry landscape is continuously evolving, changing the requirements

**Ratio of data scientists to software engineers low** → Not enough headcount with the required skill set

**New technology** → Lots of learning on every front (AI, Security, Legal, etc.,)

**Concurrent workflows, everyone is excited** → Work duplication (code and skill development)

**Roadmap realignment** → Introducing a completely new item in the roadmap

# Establish Rules on Engagement

You have a potential Use Case for LLM → follow this process

**Involve a Data Scientist in early conversations**

- Feasibility validation
- Suggesting alternative solutions (do you need an LLM?)
- What support you need
    - Prompt engineering
    - Testing plan
    - Use case - model efficacy
    - API consolidation

Presented in DataCon LA 2023 by Sanghamitra Goswami