Gen Al at Swiggy

Jairaj Sathyanarayana VP Data Science & Analytics [<u>https://linkedin.com/in/jairajs</u>]

Abhinav Nigam Director of Engineering [https://linkedin.com/in/nigamabhinav/

PROSUS AI MARKETPLACE 2023

LLMS IN PRODUCTION

A systematic approach to institutionalizing gen Al

Established a gen Al task-force (in Q1 of 2023)

Based on learnings, set up gen AI as a company-wide thrust with a single-threaded leader (starting Q2 of 2023)

- Run the thrust like any other Annual Operating Plan (AOP) program... but with some flexibility for course corrections and fuzzy bets
 - Funding-- headcount, program management \bigcirc
 - Evangelization & enablement-- hackathons, tools Ο
 - Monthly reviews with the CXO team Ο
- 3 broad classes of gen AI initiatives
 - **Discovery & Search** Ο
 - Productivity (excluding no-brainer use-cases like adopting Github Copilot)
 - Automation \bigcirc

Discovery & Search track: Catalog Enrichment Currently in production: Text generation (GPT3.5 Turbo)



- Prompt engineering + pre/post processing
- Summary/ text needs to be truthful but not overly negative

500K+ descriptions (99%+ acceptance rate); 1K+ SEO pages

	.ul 🗢 🔲		
∶afe atings)	¢	0	
itefield 🔻 ivery to Homeprestige 🔻			
Free delivery with Swiggy O	ne		
very on orders above ₹149			
25% off upto USE ONECARD100 A	5 ₹10 BOVE ₹2	2/4	
M E N U			
Search for dishes	Q	Ŷ	
galitfree Bestseller	_		
e ople are saying restaurant offers delicious di d Wheat Chapati. Overall, ta od, but prices were high.	ishes lik ste and	е	

V Pure Veg

NEW What p The food at the

Palak Paneer a quality were go

Bestsel

Want to repeat?

A/B with a few thousand restaurants under way in MUM, HYD

Discovery & Search track: Catalog Enrichment Image generation using ControlNet-guided base SD

How to generate images in the style of a given restaurant, at scale? Prompt engineering led solutions are not scalable (hard to control, not realistic, Indian context, etc.)







Background (real/generated) Image blending

Fixed text guidance for background (Background prompt: perspective view, plain background, blurred background, perspective view, professional, **Negative prompt**: 'ugly, human face, hands, limbs, unattractive, fork, spoon, chopstick, cutlery, texture, text, flowers,logo')



Item blended in the style of the given background (i.e., in the style of a given restaurant)





Discovery & Search track: Catalog Enrichment Pipeline based on custom & off the shelf modules

Text guidance for background
(Background prompt:
perspective view, plain
background, blurred
background, perspective view,
professional, Negative prompt:
'ugly, human face, hands, limbs,
unattractive, fork, spoon,
chopstick, cutlery, texture, text,
flowers,logo')

CLIP Text Encoder



Model

Inputs



LoRA checkpoints at dish-family level (for ex, pizzas, 35 such checkpoints-- 80% of order volume ~2000 images for each checkpoint training; 8 sec to

generate an image on a 16GB GPU

We get both Text2Img and Img2Img capabilities

Txt2Img

Grilled Chicken Burger













Image Blending

Name: Dry Chilli Chicken

Reference Image





Img2Img



Pipeline	Approval rate	Best-Image tag
Image Blending	51%	29%
Txt2Img	62%	26%
lmg2lmg	69%	45%

Discovery & Search track: Enhancing our Search



Discovery & Search track: Enhancing our Search The ML pipeline: fine-tuned LM based on all-mpnet-base-v2 (100M+ params)



Unsupervised fine-tuning Eg: TSDAE, MLM Noise introduced in text is reconstructed to the original; Training for a few epochs till cosine similarity on similar word pairs in our context goes up Eg: Sim("rice bath, bisibelebath") > Sim("rice bath","fried rice")

Supervised fine-tuning Eg: Contrastive learning, Multiple Negative Ranking Loss

(query, item name) pairs are used in training

Eg: South indian breakfast \rightarrow Idli, Healthy beverage \rightarrow ABC juice

Productivity track: Conversational Analytics



Rolling out to a few thousand restaurants

SIGN UP NOW!

70-90% accurate in tests; In beta with select domains-trending closer to 50%

Data/ Biz context is the long pole; Tech is doable

Productivity track: Conversational Analytics Text2Text: WhatsApp assistant for Restaurant Partners

- Restaurant owners have questions like 'How do I setup a new outlet', 'How do I mark an item out of stock', 'How do claims affect my ratings', etc.
- Navigating through a maze of dispositions on Help Center is not easy
- Conversational Q&A can help improve experience, reduce TAT and load on support teams



Productivity track: Conversational Analytics



Productivity track: Conversational Analytics Text2SQL: Intermediate-level SQL based self-serve via Slack



RAG with multiple levels of retrieval: fetch relevant tables, then relevant columns and construct the final query using meta of the fetched data

Productivity track: Conversational Analytics Text2SQL: A quick demo

Task	Accuracy
Simple logic queries	90%
Complex logic queries (multiple joins, JSON processing)	70%

enu	
ad your meta-data of columns in CSV format	Text2Sql 🐃
rag and drop file here	
Browse files	
oad your meta-data of tables description in CSV nat	
rag and drop file here imit 200MB per file - CSV	
Browse files	

×

*

Productivity track: Conversational Analytics Text2Insight: Relies heavily on codifying Biz hypotheses



Productivity track: Conversational Analytics Text2Insight: From hypotheses to answers

What are some of the reasons my order counts been *increasing* day over day?

After figuring out the intent & the search-space for approaches to answering

Order counts could increase because of higher discounts or increased ads spending, check if that has happened and call it out to the user. Also, compute order counts for each day and determine whether the orders counts on average are trending upwards. If there is a significant increase (>50% of average daily counts) for any particular days, call it out as an exceptional day in terms of order volume.

The days with order counts significantly higher than the average (more than 1.5 times the average) are:

• • •

2023-06-11: 911 orders 2023-06-18: 1078 orders 2023-06-25: 1039 orders 2023-07-02: 1012 orders

These days can be considered as exceptional in terms of order volume.

In conclusion, the increase in order counts could be due to the higher discounts given on certain days. The days with significantly higher order counts also coincide with the days with higher discounts. This suggests that discounts could be driving the increase in order counts.

Top-k embeddings (past answers)

QNLI Filtering

1.

Order counts could drop because of lowered discounts or reduced ads spending. check if that has happened and call it out the user.

2. Compute order counts for each day and determine whether the orders counts on average are trending upwards or downwards. Also call out if order volume has dropped significantly(>50% of average daily counts) for any particular days as an anomalous day in terms of order volume.

['Order counts could drop because of lowered discounts or reduced ads spending, check if that has happened and call it out the user.

"Order counts are also a function of days of the week, compare the order counts with the same days of the week in the previous week. If that's comparable, that is within 30% margin, then there is no actual drop and call that out.",

'To answer any questions regarding AOV values, using the bill_amount column to aggregate averages.',

'Drop in average order values could be due to increased average discounts that might reduce the overall order value."

'Drop in average order values could be due to increased sales proportion of a cheaper item or decreased sales proportion of an expensive item in the current duration as opposed to a few days before which can be the reference duration. Compare the sales the top selling item and their respective prices between time periods and if that indicates a drop in AOV, call that out.',

'Ads performance could see a dip in case there are other restaurants that are similar to yours in terms of cuisines offered that have increased their ads spending in recent time.',

'Compute order counts for each day and determine whether the orders counts on average are trending upwards or downwards. Also call out if order volume has dropped significantly(>50% of average daily counts) for any particular days as an anomalous day in terms of order volume.']

Automation track

Automation of dispositions that do NOT require a lot of context

Email res

				Explore	your plan
mail settings Q. Search email address			@ A1	dvanced settings New support e	Email settings mail You can manage all your support mailboxes here. Any email sent to your support email address for example: hello@acme.com gets automatically converted into a ticket that you can get working on.
Product	Name	Email address	Group	Status	To create a new support email box, click
🛨 Global Support Mail	Hello Acme	hello@acme.com Gmail Reauthorization required. Click here	-	Verify	"New support email" under email setting
Global Support Mail	Request Acme	☑ requesteracme@acme.com		Verify	Help me set up my support email
		Gmail Reauthorization required. Click here			Using multiple mailboxes
					You can add unlimited incoming and outgoing mailboxes (like info@yourcompany.com, sales@yourcompany.com, etc.) in your helpdesk. You can even setup Freshdesk Support Desk to automatically assign emails from each mailbox to a specific group.
					Reply emails from Freshdesk Support De will automatically use your respective support email IDs as the From email address. Creating an SPF record in your

execute SOP, generate empathetic response

Automation track: LLM customer-service bot

7:47 \ mb ♥ •	7:56 🂭 🔪 mater	
Chat with Swiggy	← Chat with Swiggy	
I have a coupon related query for non food	You are now chatting v	vith santosh.rani
issues 07:35 PM	My name is Santosh rani, you today. Give me a mo is the next course of acti	and I'll be assistir ment to check wh on.
Hi there! I will be assisting you with this. Please let us know the issue you are facing?	0	Sep 15, 04:54
I am unable to find my coupon	Delivery guy came back has bike issue	with the order as
My coupon did not work as expected	Thank you for waiting!	Sep 10, 04
I forgot to apply my coupon	9 Sep 15, 04:55 PM	
I am unable to apply my coupon	We are sorry about the tr	ouble faced. pp 15, 04:56 PM
I want to know the details of my coupon	I tried contacting the del	very partner, but
07:35 PM		Sep 15, 04:57
	Could you please confirm partner returned your ord or not?	n if the delivery der pickup locatio
	0	Sep 15, 04:57
oh! Looks like this conversation is no longer valid.	As you seem to have ste chat, I am going to conta registered number.	oped away from tł ct you on your
ii you still need assistance on this query	•	Sep 15, 04:58
CHAT WITH US	This has been marked as c other queries regar	losed. If you have ding this issue
		<

Decision-tree based bot

Human agent conversation

LLM powered bot

Automation track: LLM customer-service bot

Key challenges

- Dynamic data
- Steerability/driving towards closure
- Jailbreaking/hallucinations are expensive
- Free form text related issues
 - Out of order input
 - Conversational memory
 - Context shift

Automation track: LLM customer-service bot Key learnings

- Data intensive use cases
 - Latency + token costs _
 - Remedies _
 - Context compression using LLM
 - Vectorization/RAG (Retrieval Augmented Generation) pipeline _
 - Streaming API for gathering LLM responses
- Jailbreaking
 - System message to limit the context -
 - Moving system message to the end
 - Few shot prompting _
- Deploying for customers
 - VPC with firewalls and DDoS protection -
 - Tenant isolation

Some learnings so far

- It took us about 3-4 months of wandering to land on potentially high ROI items. Until then we ended up flailing quite a bit
 - Conserve your bandwidth-- limit how many inbounds you entertain (both internal and external) \bigcirc
 - There's a delta in effort/quality between a cool hackathon project/ external demo and it's production version. Ο Set/ manage expectations

- For cases which are not real-time, hard to beat GPT out of the box Cost <> ROI trade-off 0
 - Hallucination was a real problem; Lots of internal-user testing & guardrailing \bigcirc
 - Using GPT directly from OpenAI quickly led to governance difficulties. Moved to Azure AI \bigcirc
 - This also means fine-tuning is usually a last resort/ expensive choice Ο

• We didn't see a lot of pull for customer-facing conversational interfaces ('chatbot for food ordering / table reservations')

Some learnings so far

Be dispassionate and replace if something better comes along Replaced Image Blending with Img2Img despite having spent a decent amount of time on the former

- Data and metadata are paramount. Especially for conversational analytics use-cases, we had to source and codify a lot of context / knowledge sitting in Business/ Ops folks' heads
 - 'Copilot for X' use-cases are tempting but quickly run into data difficulties Ο

- Be pragmatic and be patient
 - Not easy to find/ land slam-dunk use-cases
 - Relatively easier to find automation tasks (i.e., efficiency improvement/ cost reduction) vs. top-line ones Ο
 - Human-in-the-loop/ Ops validation is almost always necessary \bigcirc

THANK YOU!

Don't forget to: Swiggy today! Get your Swiggy One!

PROSUS AI MARKETPLACE 2023

LLMS IN PRODUCTION