

Generative Models For Data Augmentation

Brandon Trabucco, 4/9/25



Machine
Learning
Department

Carnegie Mellon University
School of Computer Science

Towards Generative Data Augmentation



Horizontal Flip



RandAugment
(Cubuk et al., 2019)



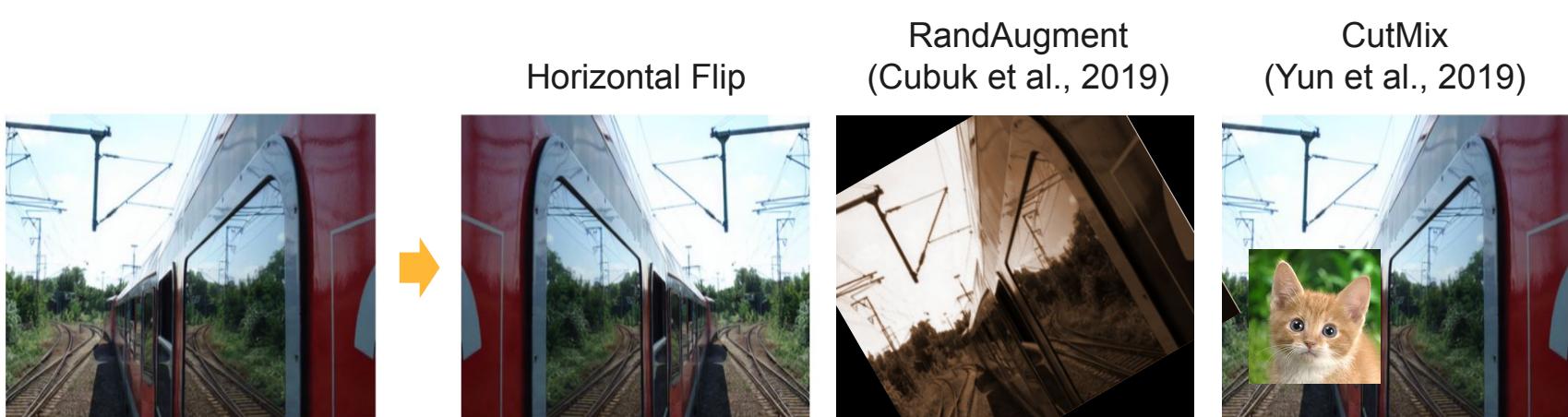
CutMix
(Yun et al., 2019)



[1] Cubuk et al., RandAugment: Practical automated data augmentation with a reduced search space, NeurIPS 2020.

[2] Yun et al., CutMix: Regularization Strategy to Train Strong Classifiers with Localizable Features, ICCV 2020.

Towards Generative Data Augmentation



- Data augmentation often **requires a good intuition** about your dataset.

[1] Cubuk et al., RandAugment: Practical automated data augmentation with a reduced search space, NeurIPS 2020.

[2] Yun et al., CutMix: Regularization Strategy to Train Strong Classifiers with Localizable Features, ICCV 2020.

We Need Augmentations That **Adapt To Your Dataset**



[1] Cubuk et al., RandAugment: Practical automated data augmentation with a reduced search space, NeurIPS 2020.

[2] Yun et al., CutMix: Regularization Strategy to Train Strong Classifiers with Localizable Features, ICCV 2020.

We Need Augmentations That **Adapt To Your Dataset**



[1] Cubuk et al., RandAugment: Practical automated data augmentation with a reduced search space, NeurIPS 2020.

[2] Yun et al., CutMix: Regularization Strategy to Train Strong Classifiers with Localizable Features, ICCV 2020.

Are Generative Models Capable Enough?



GAN, 2014 [3]



DCGAN, 2016 [4]



BigGAN, 2019 [5]



StableDiffusion, 2022 [6]

Nightmares before final exams

Works of art

- Generative models have developed **astounding levels** of **photo-realism**.

[3] Goodfellow, Ian, et al., Generative Adversarial Networks, NeurIPS 2014.

[4] Radford, Alec, et al., Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks, ICLR 2016.

[5] Brock, Andrew, et al., Large Scale GAN Training for High Fidelity Natural Image Synthesis, ICLR 2019.

[6] Rombach, Robin, et al., High-Resolution Image Synthesis with Latent Diffusion Models, CVPR 2022.

Can We Harness Image-Editing Techniques?

Why is this a good idea?

- **Data scarcity**: we can sample as many images as we need.
- **Semantics**: we can choose what to edit.

Key Idea: target and modify key semantic properties of images.



- Generations from a recent image-editing technique that respects high-level structure [7].

[7] Hertz, Amir, et al., Prompt-to-Prompt Image Editing with Cross Attention Control, arXiv 2022.

Generative Models Capture New Invariances



Rotate



Flip + Saturate



Standard
Augmentation



Lemon → Apple



Lemon → Pistachio

“Data Augmentation”
With Generative Models

- Generative models can target **semantic properties** of real images that we expect and desire models to be invariant to.

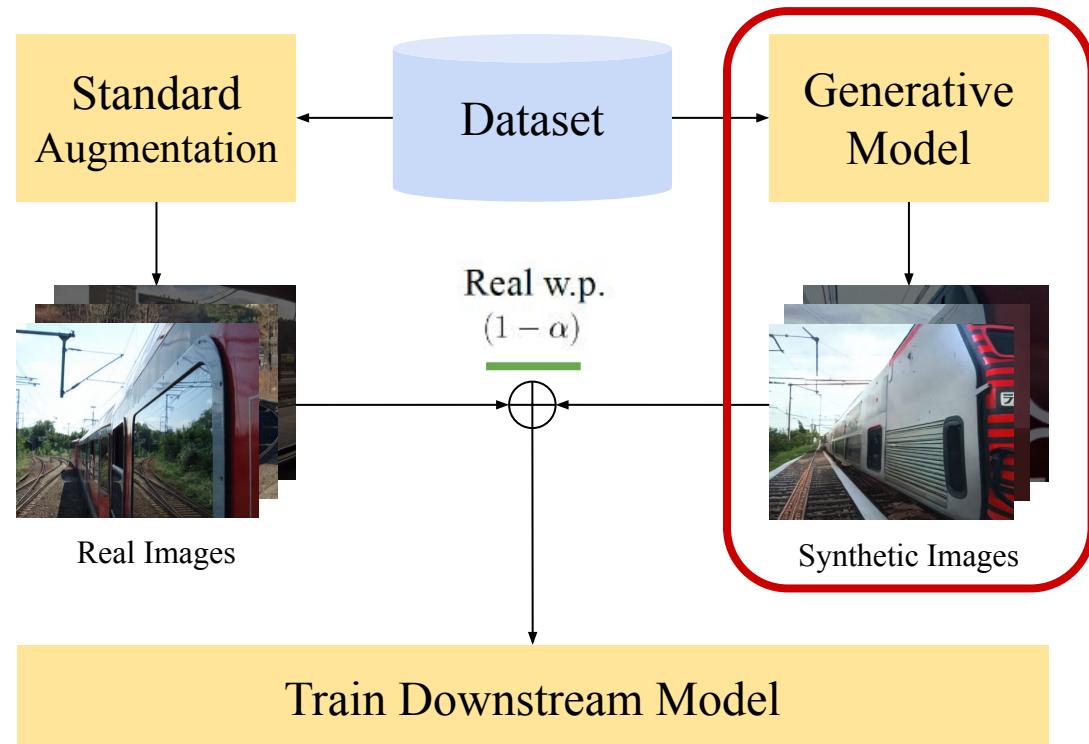
[7] Hertz, Amir, et al., Prompt-to-Prompt Image Editing with Cross Attention Control, arXiv 2022.

The Generative Data Augmentation Pipeline

- Mix real and generated images, weighted by an **augmentation probability**.

New Axes To Explore:

- What image-editing technique shall we use?
- How much structure do we preserve in the data?



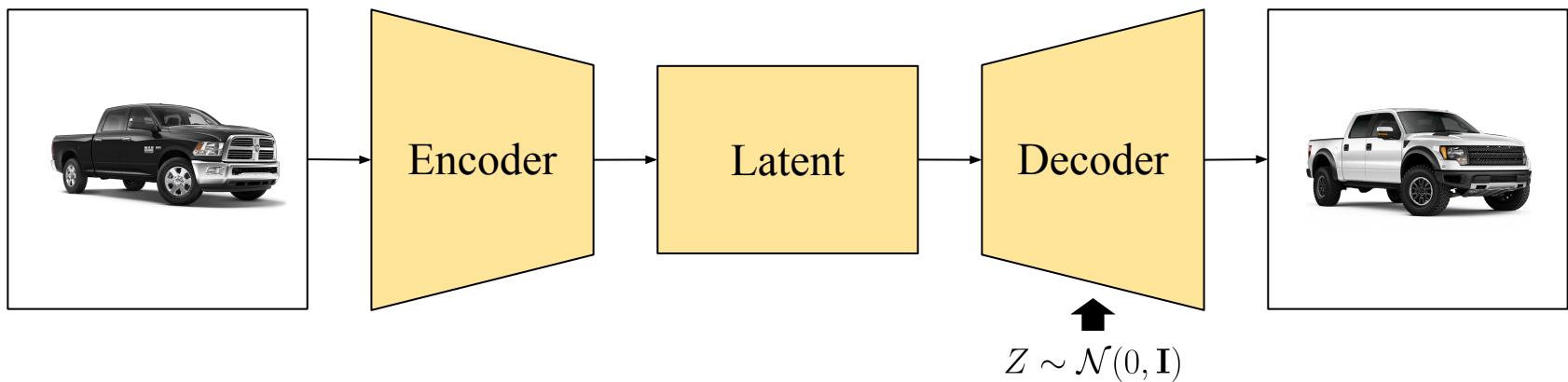
How Do We **Control** The Edited Property?

- Modern generative models have prompts, but this is a **recent development**.
- Before such models, researchers had explored methods to **infer the distribution of edits** purely from images.



How Do We Control The Edited Property?

- How can we **infer the distribution of edits** purely from images.
- **DAGAN:** model this as an auto-encoding problem.

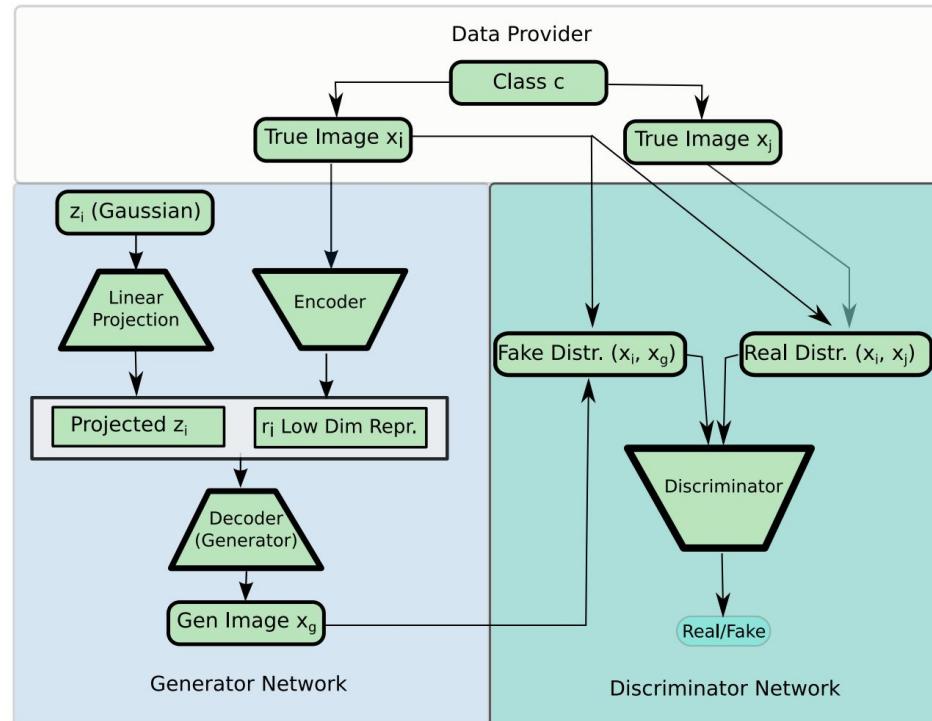


Data Augmentation Generative Adversarial Networks

Modelling Task:

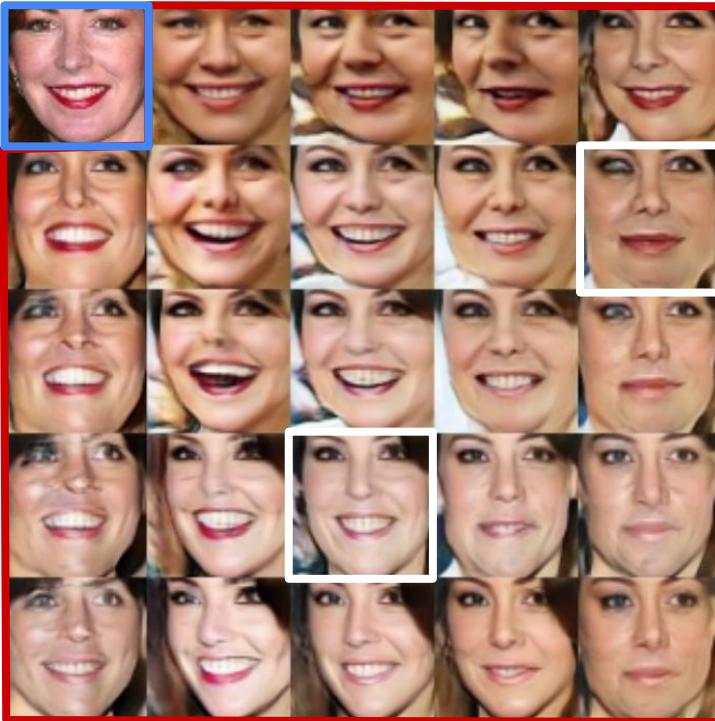
Generate an image from the same class as image X_i , *without observing the class*

- Generator implemented as a **UNet** mapping X_i to a generation X_g .
- Discriminator predicts if X_i, X_g are the same class.



DAGAN: Example Generations & Results

Real



Fake

Face DAGAN Augmented Classification		
Experiment ID	Samples Per Class	Test Accuracy
VGG-Face_Standard	5	0.0446948
VGG-Face_DAGAN_Augmented	5	0.125969
VGG-Face_Standard	15	0.39329
VGG-Face_DAGAN_Augmented	15	0.429385
VGG-Face_Standard	25	0.579942
VGG-Face_DAGAN_Augmented	25	0.584666

- Consistent improvement when tested on held-out classes.
- Diminishing gains when many samples per class are available.
- Generations are **inflexible**.

How **Flexible** Is Data Augmentation GAN?



Rotate



Flip + Saturate



DAGAN

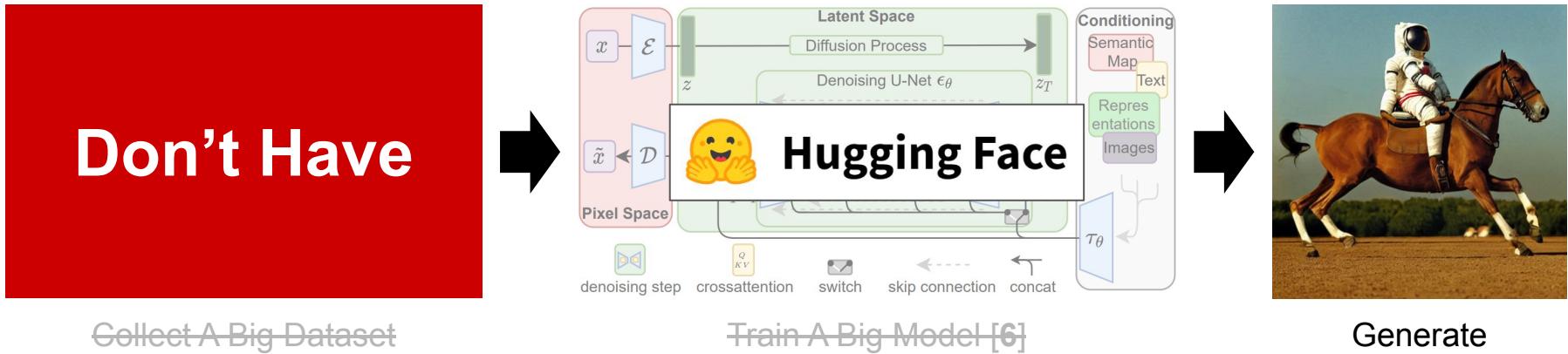
$Z \sim \mathcal{N}(0, \mathbf{I})$



- **Controllable** extent and randomness.
- Highly modular.
- Generates images **all at once** with little control over the layout and content.
- Requires **training** a GAN.

DA-GAN Questions?

Can We Avoid Training A New Generative Model?



- Let's use **pre-trained** image generative models for data augmentation.
- Several powerful models: Imagen, GLIDE, **Stable Diffusion**.

[6] Rombach, Robin, et al., High-Resolution Image Synthesis with Latent Diffusion Models, CVPR 2022.

DA-Fusion: Data Augmentation With Diffusion

- **Key idea:** **shared context** in your images controls the augmentation.

[3] Rombach et al., High-Resolution Image Synthesis with Latent Diffusion Models, CVPR 2022.

[7] Rinon, Gal, et al., An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion, CVPR 2022.

DA-Fusion: Data Augmentation With Diffusion

- **Key idea:** **shared context** in your images controls the augmentation.

$$\max_y \log P(\text{cat with bowtie} | \text{a photo of a } y)$$

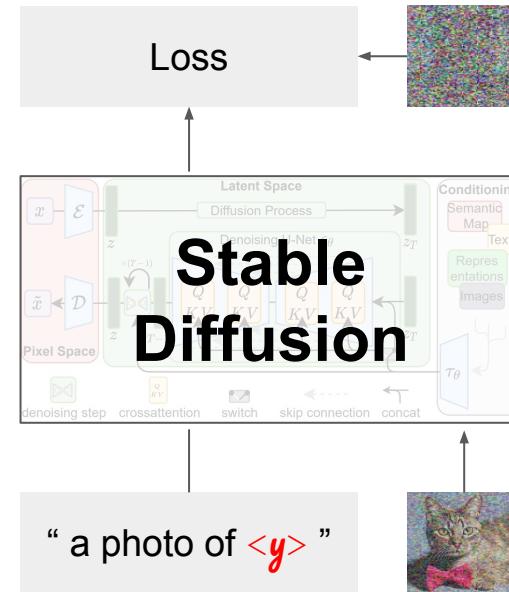
[3] Rombach et al., High-Resolution Image Synthesis with Latent Diffusion Models, CVPR 2022.

[7] Rinon, Gal, et al., An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion, CVPR 2022.

DA-Fusion: Data Augmentation With Diffusion

- **Key idea:** **shared context** in your images controls the augmentation.

$$\max_y \log P(\text{cat with bowtie} | \text{a photo of } y)$$



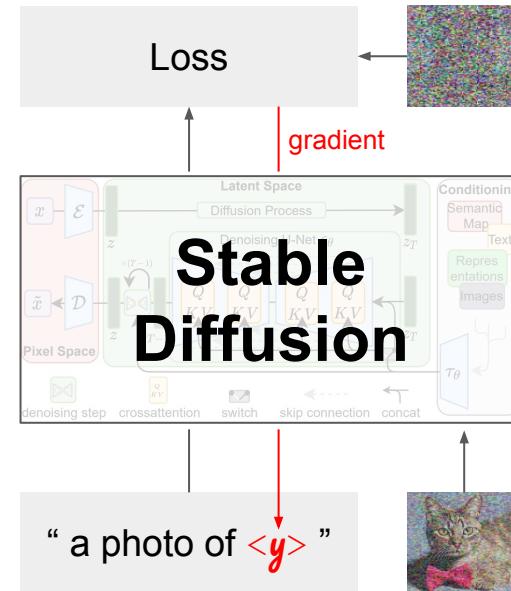
[3] Rombach et al., High-Resolution Image Synthesis with Latent Diffusion Models, CVPR 2022.

[7] Rinon, Gal, et al., An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion, CVPR 2022.

DA-Fusion: Data Augmentation With Diffusion

- **Key idea:** **shared context** in your images controls the augmentation.

$$\max_y \log P(\text{cat with bowtie} | \text{a photo of } y)$$



$\langle y \rangle \sim \text{cat wearing a red bow-tie}$

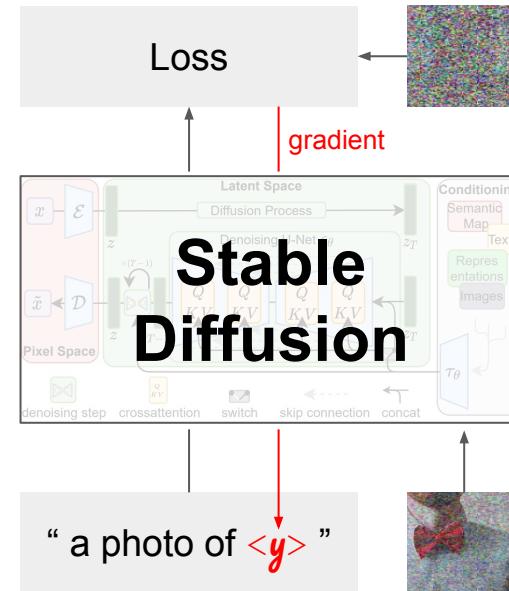
[3] Rombach et al., High-Resolution Image Synthesis with Latent Diffusion Models, CVPR 2022.

[7] Rinon, Gal, et al., An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion, CVPR 2022.

DA-Fusion: Data Augmentation With Diffusion

- **Key idea:** **shared context** in your images controls the augmentation.


$$\max_y \log P(\text{a photo of a } y) | \text{ a photo of a } y$$



$\langle y \rangle \sim \text{wearing a red bow-tie}$

[3] Rombach et al., High-Resolution Image Synthesis with Latent Diffusion Models, CVPR 2022.

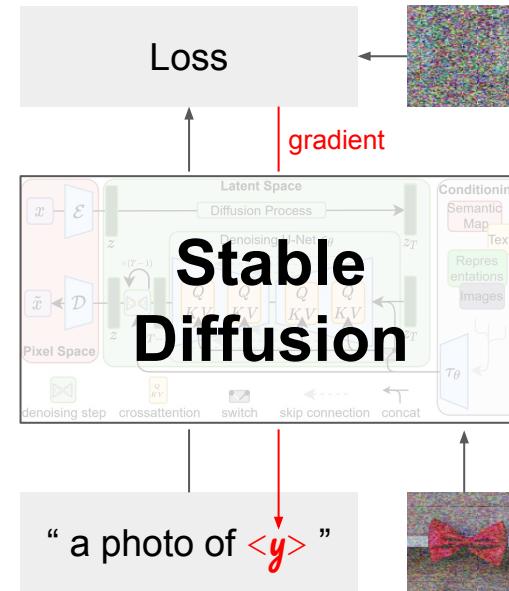
[7] Rinon, Gal, et al., An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion, CVPR 2022.

DA-Fusion: Data Augmentation With Diffusion

- **Key idea:** **shared context** in your images controls the augmentation.



$$\max_y \log P(\text{red bow-tie} | \text{a photo of } y)$$



$\langle y \rangle \sim \text{red bow-tie}$

[3] Rombach et al., High-Resolution Image Synthesis with Latent Diffusion Models, CVPR 2022.

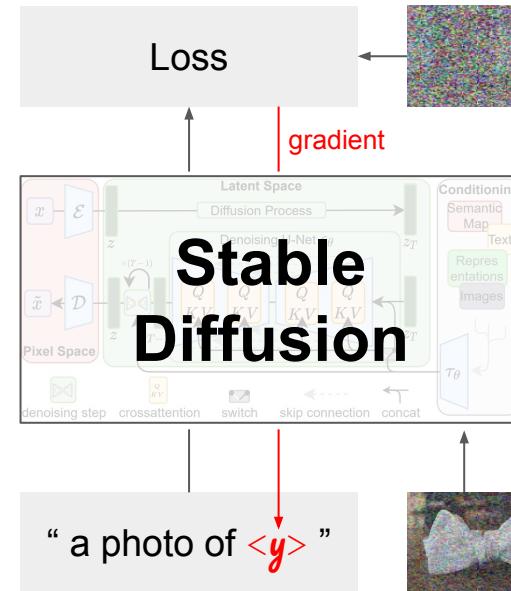
[7] Rinon, Gal, et al., An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion, CVPR 2022.

DA-Fusion: Data Augmentation With Diffusion

- **Key idea:** **shared context** in your images controls the augmentation.



$$\max_y \log P(\text{bow-tie} | \text{a photo of } y)$$



$\langle y \rangle \sim \text{bow-tie}$

[3] Rombach et al., High-Resolution Image Synthesis with Latent Diffusion Models, CVPR 2022.

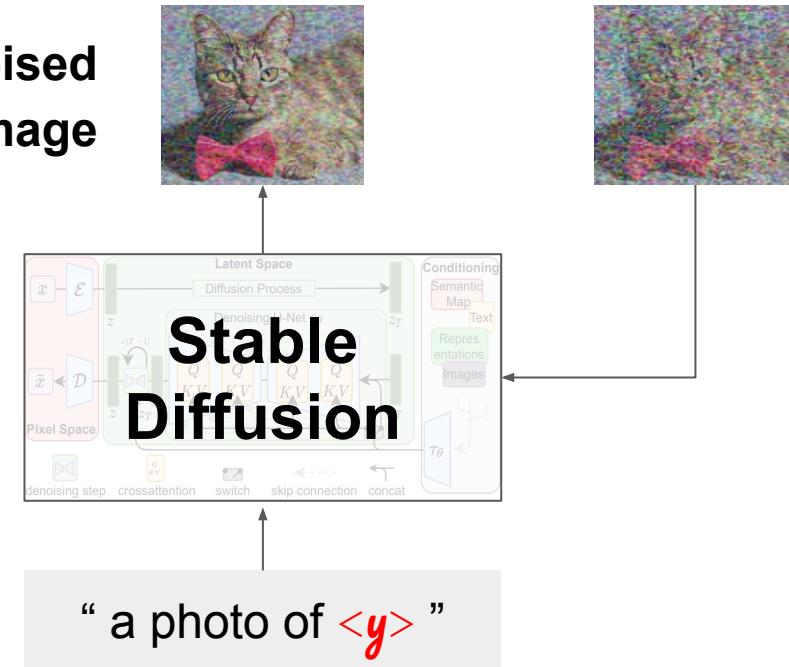
[7] Rinon, Gal, et al., An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion, CVPR 2022.

DA-Fusion: Data Augmentation With Diffusion

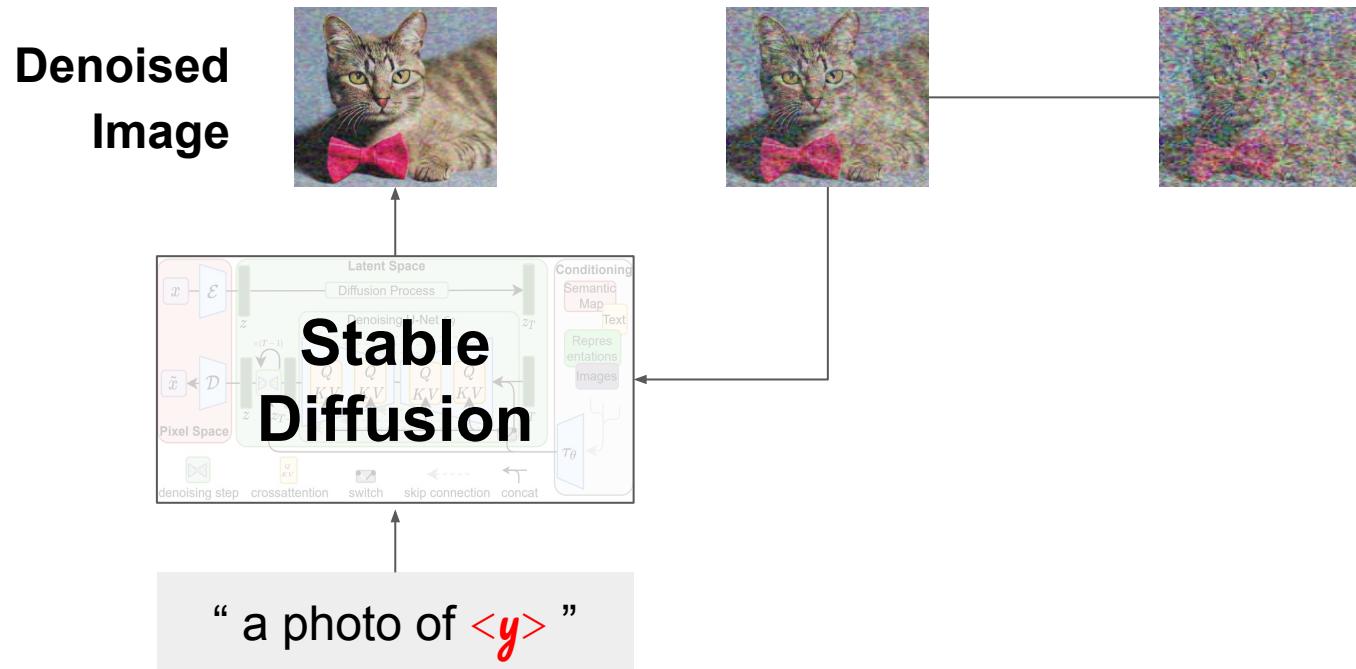
Source
Image



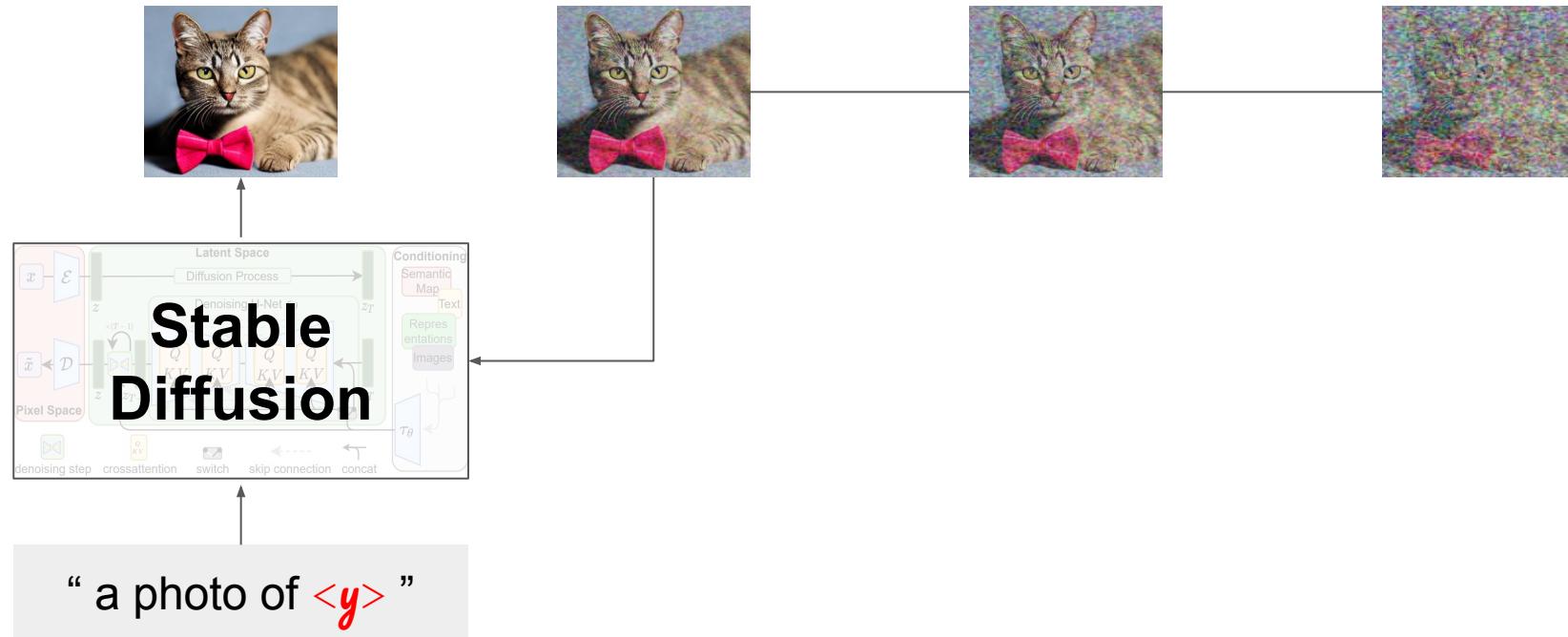
DA-Fusion: Data Augmentation With Diffusion



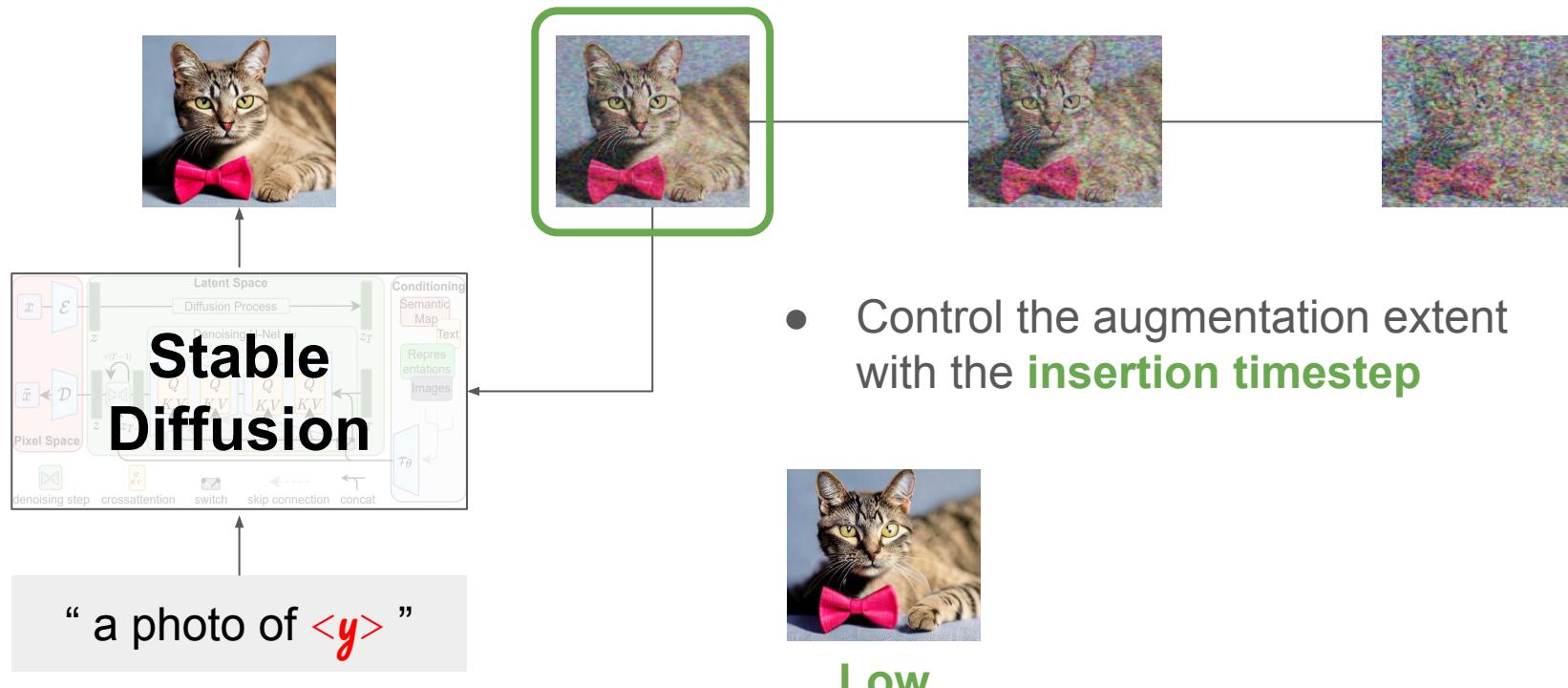
DA-Fusion: Data Augmentation With Diffusion



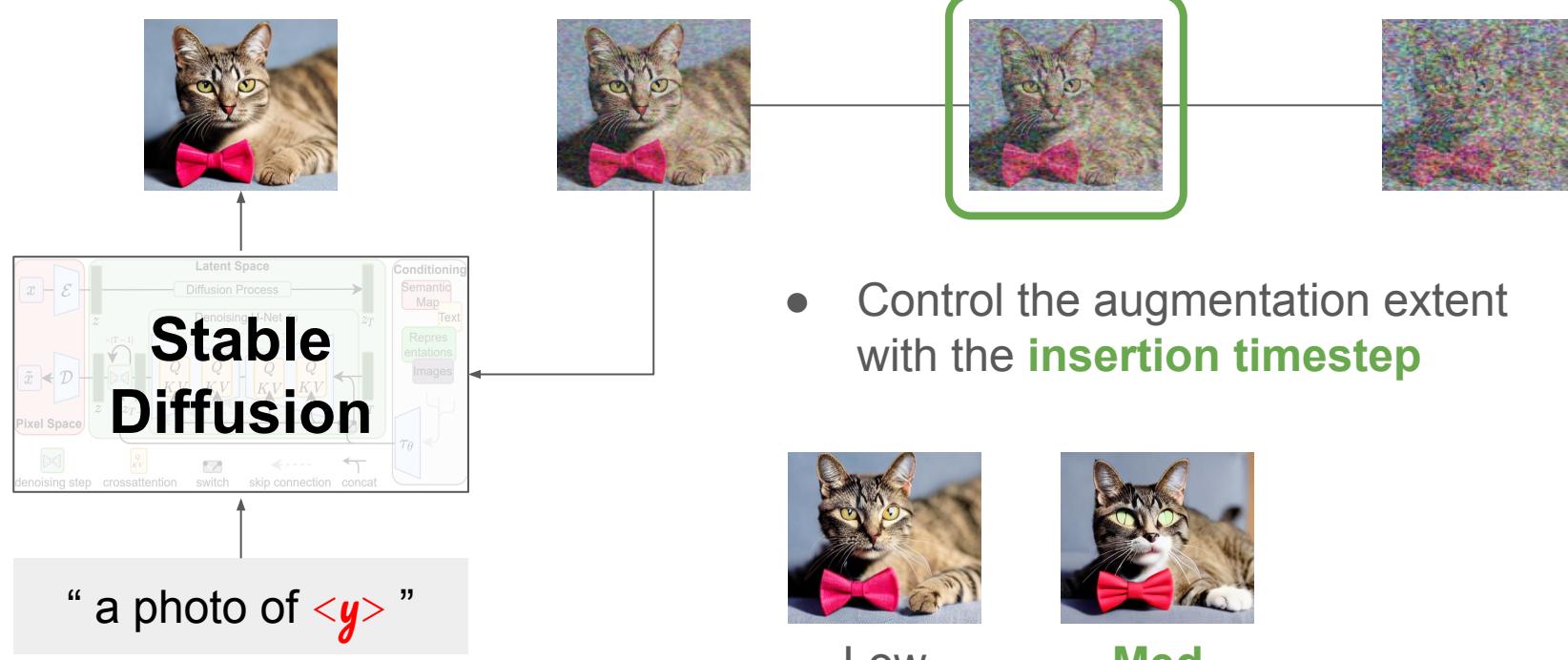
DA-Fusion: Data Augmentation With Diffusion



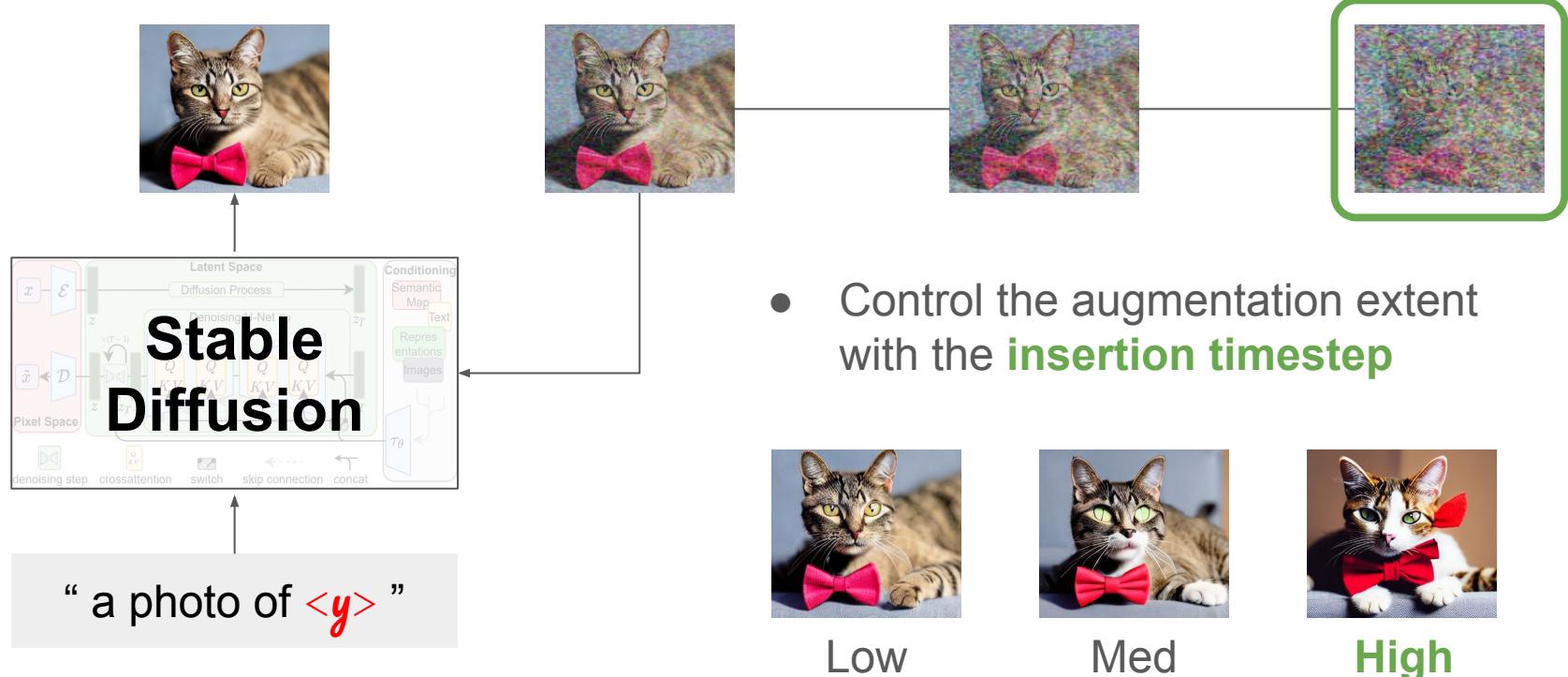
DA-Fusion: Data Augmentation With Diffusion



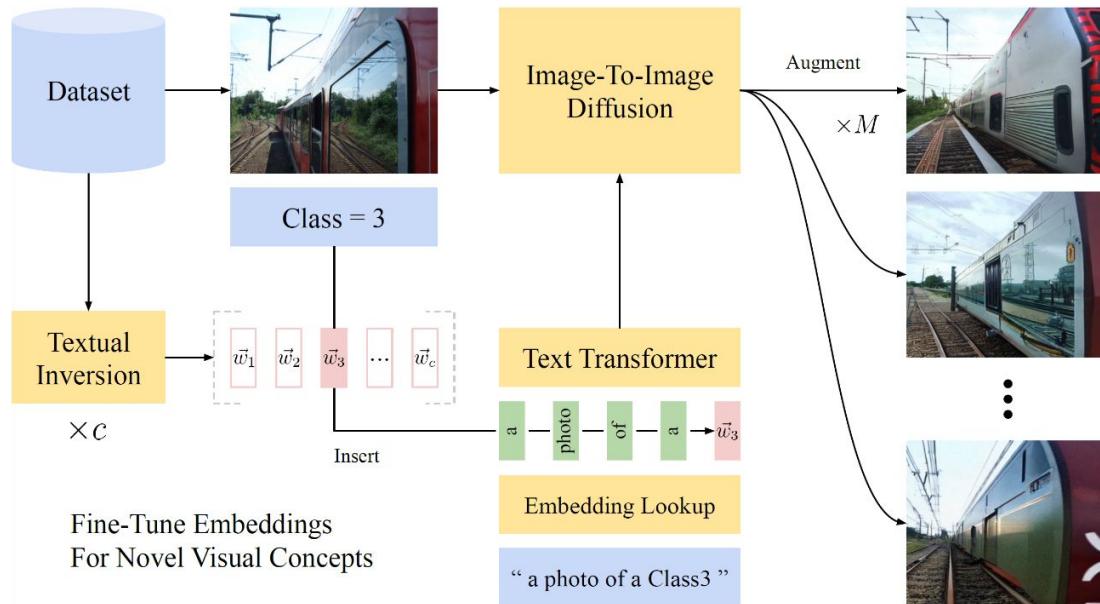
DA-Fusion: Data Augmentation With Diffusion



DA-Fusion: Data Augmentation With Diffusion

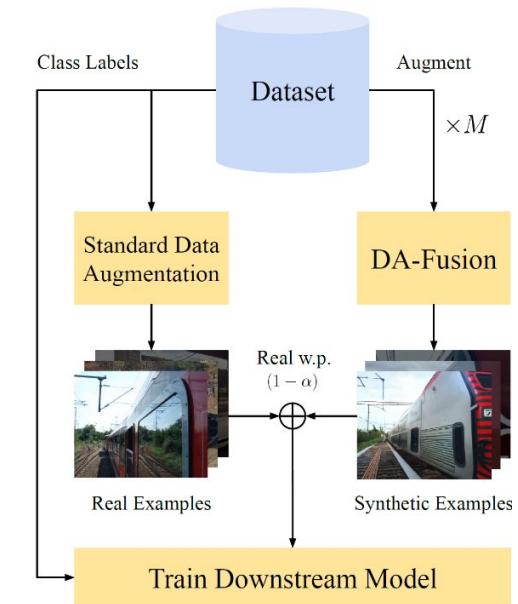


DA-Fusion: Data Augmentation With Diffusion



Fine-Tune Embeddings
For Novel Visual Concepts

- [7] Rinon, Gal, et al., An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion, CVPR 2022.
- [8] Chenlin, Meng, et al., SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations, ICLR 2022.





DA-Fusion Questions?

How Do We Evaluate DA-Fusion?

Question: How much do augmentations from DA-Fusion improve classification?

How Do We Evaluate DA-Fusion?

Question: How much do augmentations from DA-Fusion improve classification?

- Six few-shot classification tasks from literature



Common Concepts



Fine-Grain Concepts

How Do We Evaluate DA-Fusion?

Question: How much do augmentations from DA-Fusion improve classification?

- Six few-shot classification tasks from literature and **one we contribute**.



Common Concepts



Fine-Grain Concepts



Novel Concepts

How Do We Evaluate DA-Fusion?

Question: How much do augmentations from DA-Fusion improve classification?

- Given a handful of **real images**

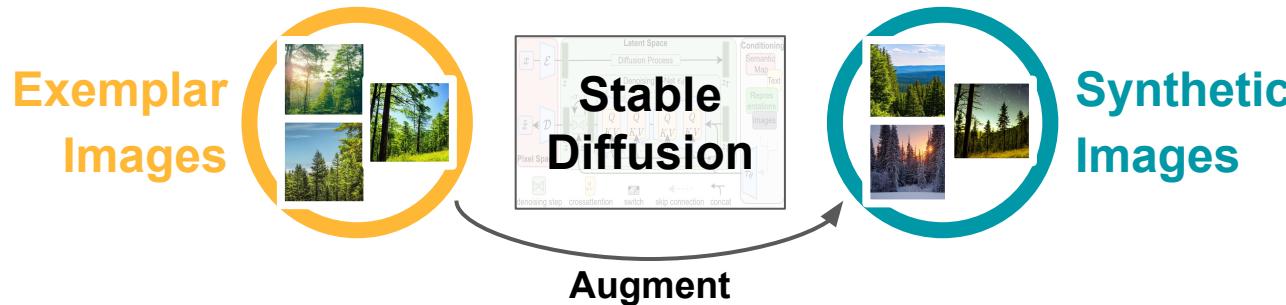
Exemplar
Images



How Do We Evaluate DA-Fusion?

Question: How much do augmentations from DA-Fusion improve classification?

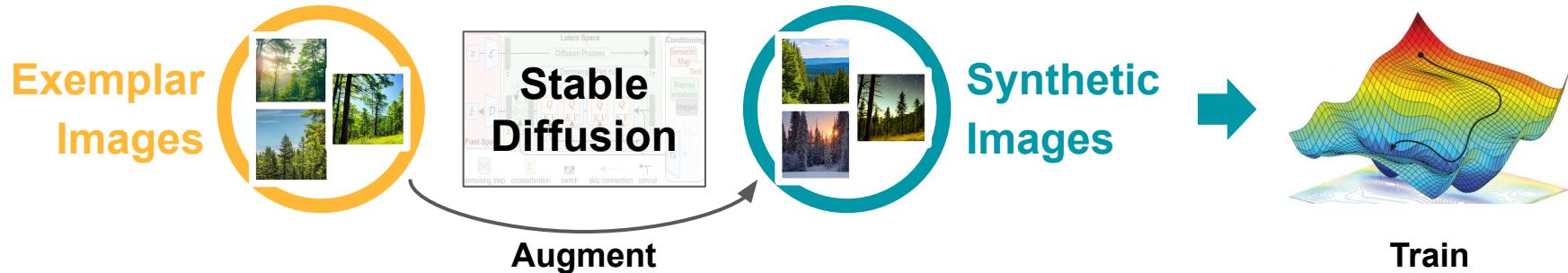
- Given a handful of **real images**, generate **augmentations**



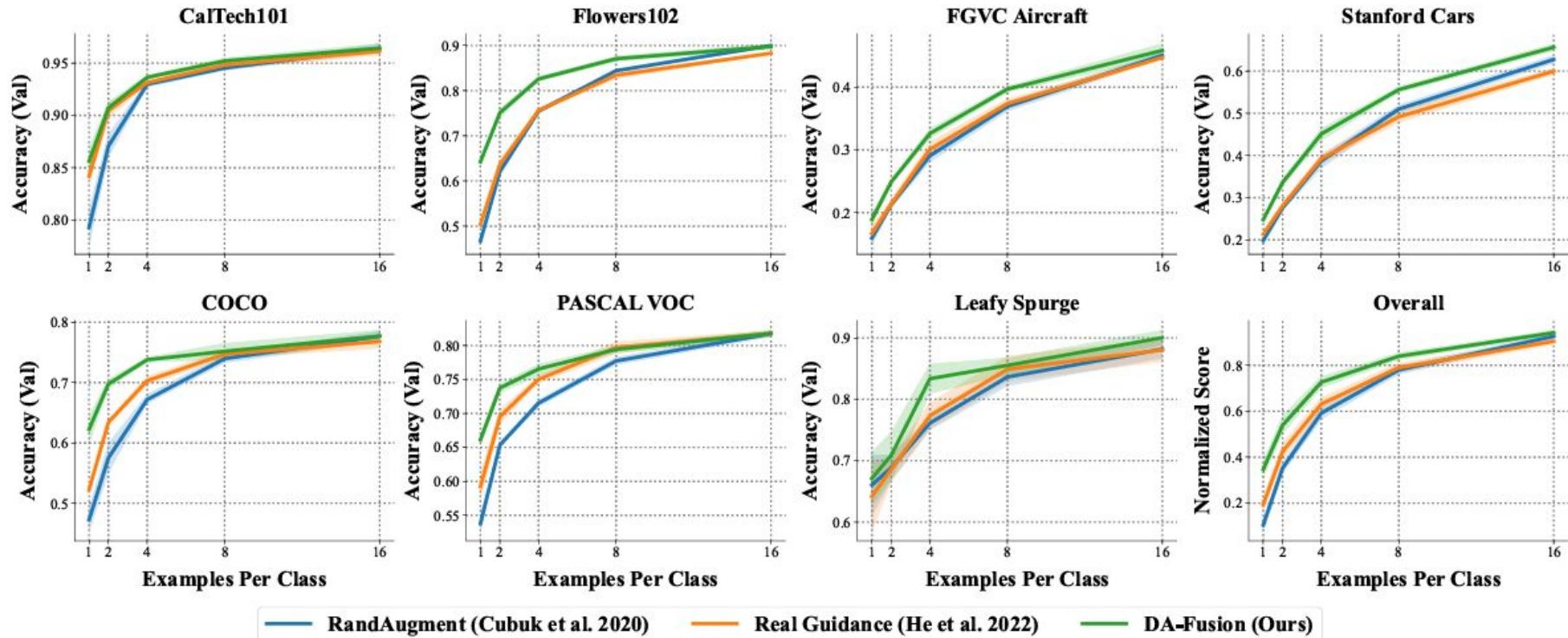
How Do We Evaluate DA-Fusion?

Question: How much do augmentations from DA-Fusion improve classification?

- Given a handful of **real images**, generate **augmentations**
- Train classifiers on a mix of real and synthetic data



DA-Fusion Improves Few-Shot Learning



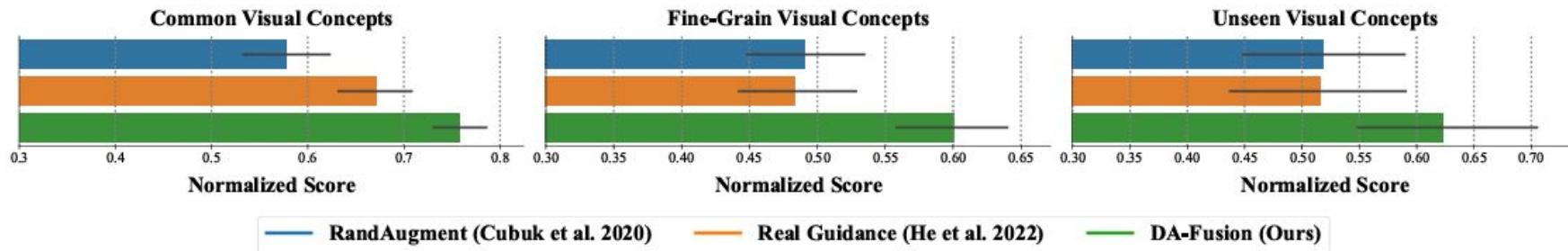
[1] Cubuk et al., RandAugment: Practical automated data augmentation with a reduced search space, NeurIPS 2020.

[2] He et al., Is synthetic data from generative models ready for image recognition?, ICLR 2023.

DA-Fusion Has Consistent Performance



Leafy
Spurge



- DA-Fusion has **strong performance** for **all types of concepts**.

[1] Cubuk et al., RandAugment: Practical automated data augmentation with a reduced search space, NeurIPS 2020.

[2] He et al., Is synthetic data from generative models ready for image recognition?, ICLR 2023.

Strong Performance ✓

Strong Performance ✓

How Do You Control The Augmentation?

When Is Additional Control Necessary?

Real images are often cluttered with **distracting concepts**.

$$\max_y \log P(\text{dog and cat} \mid \text{a photo of a } y)$$



When Is Additional Control Necessary?

Real images are often cluttered with **distracting concepts**.

$$\max_y \log P(\text{dog and cat} | \text{a photo of a } y)$$



Which concept should DA-Fusion generate: cats and/or dogs?

How Do You Control What Concept Is Learned?

Implicit Solution: **select better images** without **distracting concepts**.

$$\max_y \log P(\text{dog and cat} | \text{a photo of a } y)$$



How Do You Control What Concept Is Learned?

~~Implicit Solution: select better images without distracting concepts.~~

$$\max_y \log P(\text{dog and cat} | \text{a photo of a } y)$$



This might be **costly**, what else can we do?

How Do You Control What Concept Is Learned?

Explicit Solution: **prompt with context** about the objects you want ignored.



How Do You Control What Concept Is Learned?

Explicit Solution: **prompt with context** about the objects you want ignored.

- **Why?** Prompts can supplement information **when images have ambiguity**.



How Do You Control What Concept Is Learned?

Explicit Solution: **prompt with context** about the objects you want ignored.

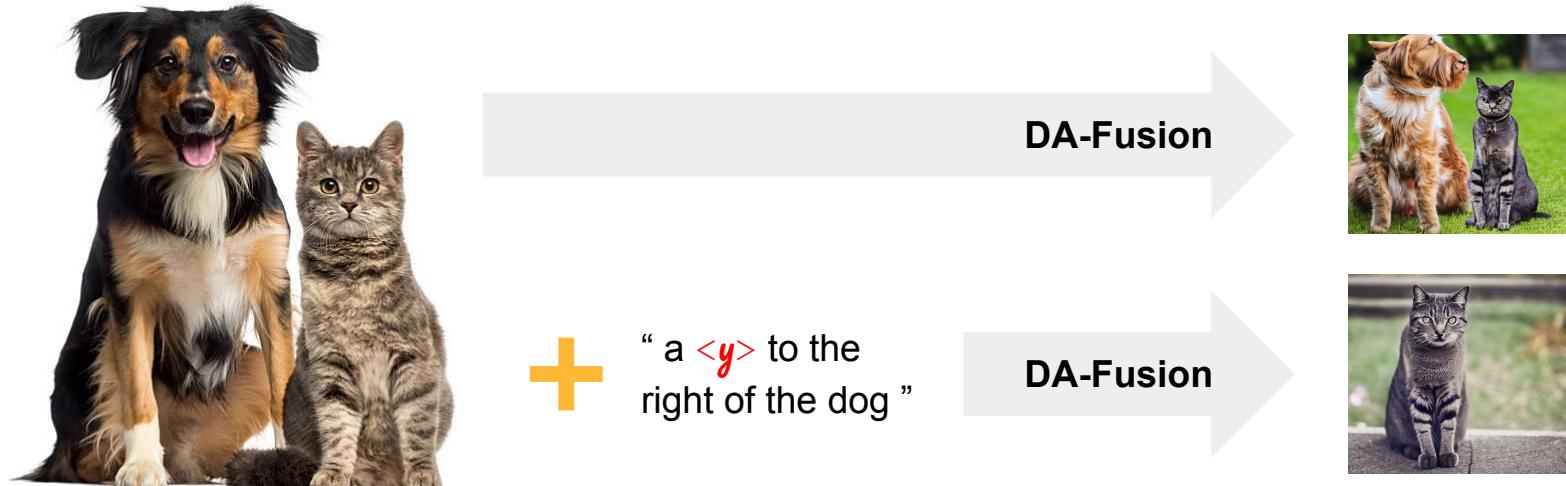
- **Why?** Prompts can supplement information **when images have ambiguity**.



How Do You Control What Concept Is Learned?

Explicit Solution: **prompt with context** about the objects you want ignored.

- **Why?** Prompts can supplement information **when images have ambiguity**.



Results Questions?

Leafy spurge (Euphorbia esula): A Problematic Weed in N. America

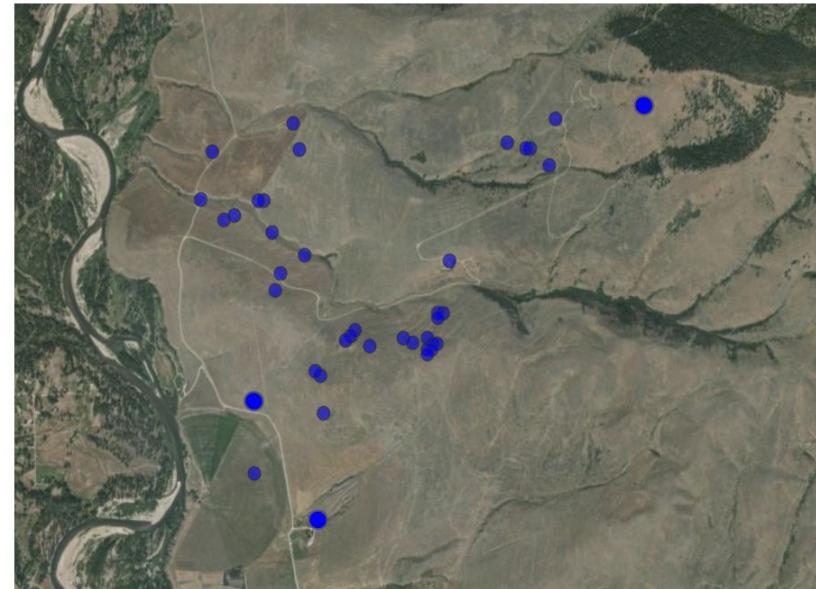


Photo credit: Montana State University



Locations of spurge surveys in western Montana

- We surveyed 40 sites varied in land-use history and plant community composition



Botanists verified spurge presence at survey sites

- We searched for spurge within nine 10 x 10 meter plots

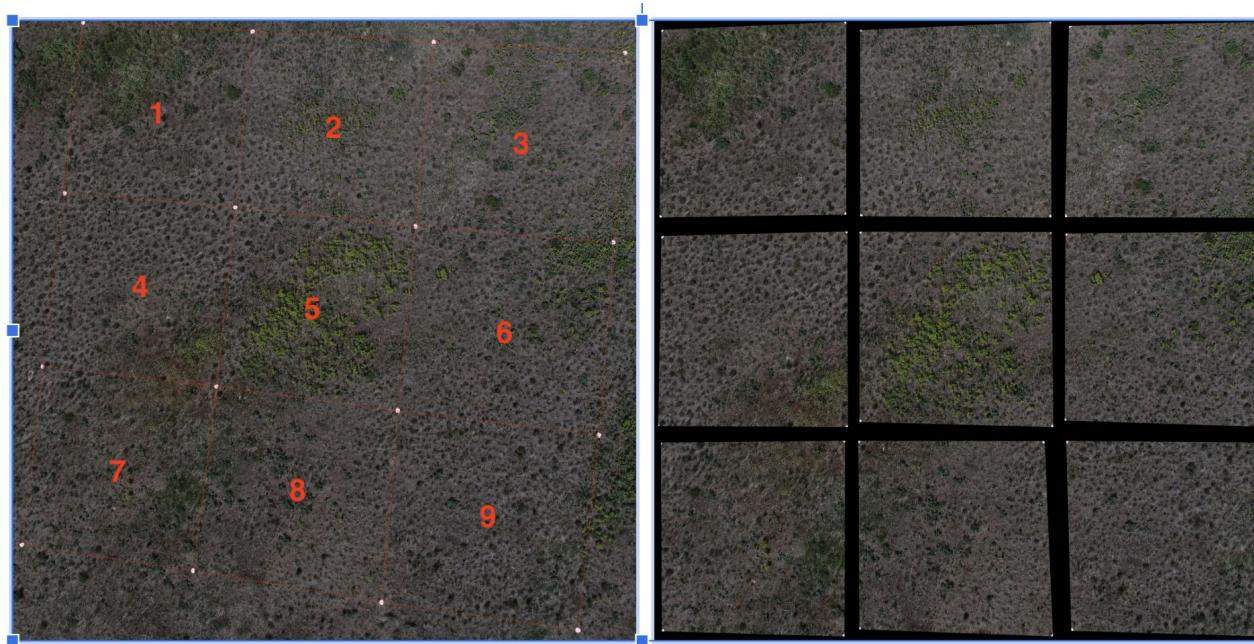


Drone then imaged surveyed areas



Post-processing: Identify plot boundaries

- Example: surveyed nine plots at each site. Markers visible in the image were used to crop plots.



Post-processing: Four-crop and verify

- We further subdivided imagery into quarters of 250×250 pixels in size (approximately 3.5×3.5 m).

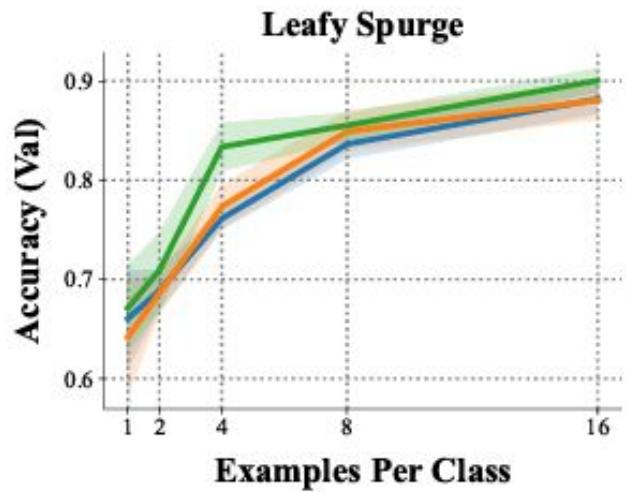


Aerial spurge images distinct from LAION dataset

- Top-down imagery of prairie ecosystems from 50m above the ground.
- Existing outside the domain of LAION and other foundation model training sets. ‘
- This makes these data well-suited for few-shot research.



Results + Synthetic Generations



- RandAugment (Cubuk et al. 2020)
- Real Guidance (He et al. 2022)
- DA-Fusion (Ours)



original generation



+ mountain



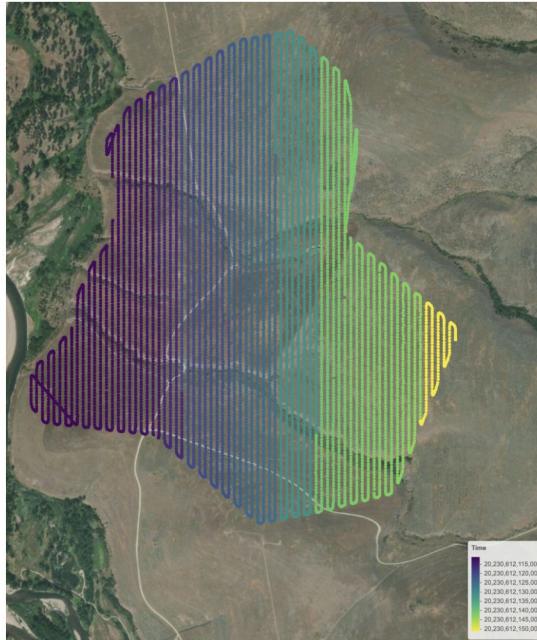
+ grassland



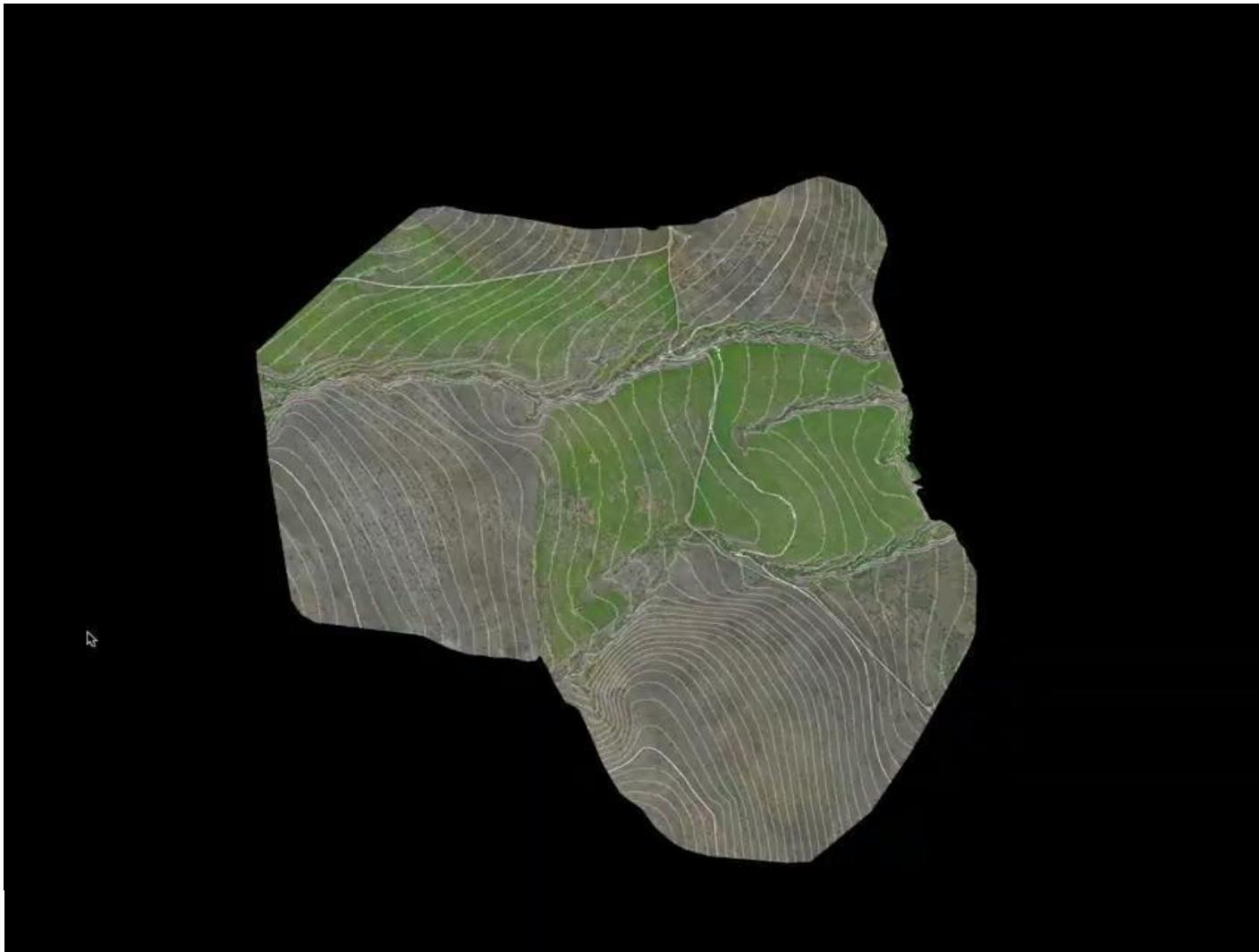
+ snow

Leafy Spurge Dataset V2

- We are working on the next version of the spurge dataset that is well-suited for mapping spurge presence at scale with a public release.







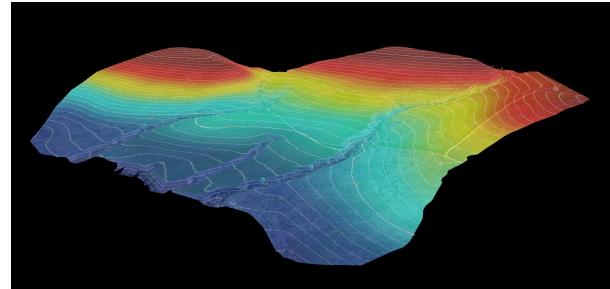


Fine-scale Classification Maps

- Fine-scale classification will enable more effective management of leafy spurge and rapidly respond to these outbreaks.
- We hope to extend this approach to other species to monitor ecosystems at scale and rapidly respond to ecological change in near real-time.



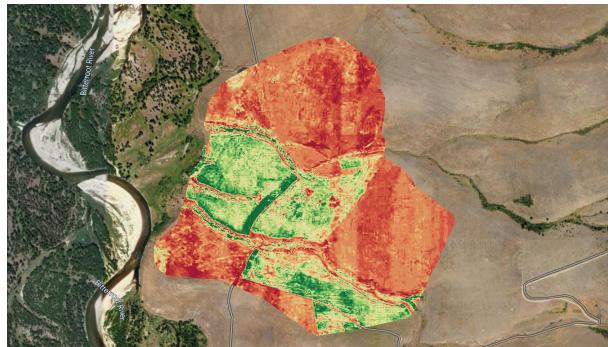
Ecologist Collaborators Offer a Rich Source of Unique Datasets



Elevation paired with RGB



Bear individual retrieval



Non-visible spectra (plant health)



Landscape change detection

Leafy Spurge Questions?

Can We Also Obtain Perceptual Annotations?



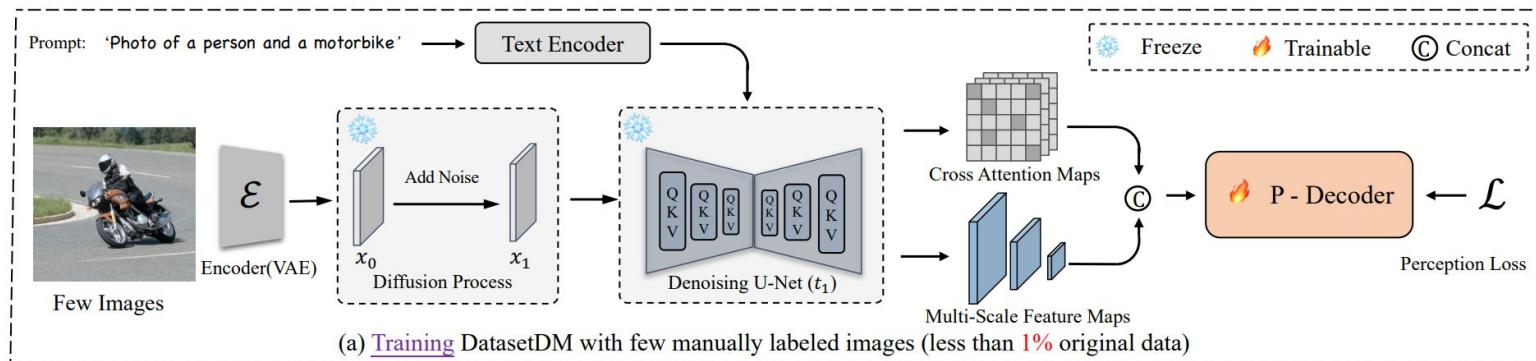
- In many data-limited settings, pixelwise labels matter, not just the class.
- Do generative models have an **inner representation** for these?
- And, can we **generate annotations** alongside the image?

[8] Weijia Wu, et al., DatasetDM: Synthesizing Data with Perception Annotations Using Diffusion Models, NeurIPS 2023.

DatasetDM: Generating Images And Annotations

Key idea: the Diffusion Model has great features, just use these.

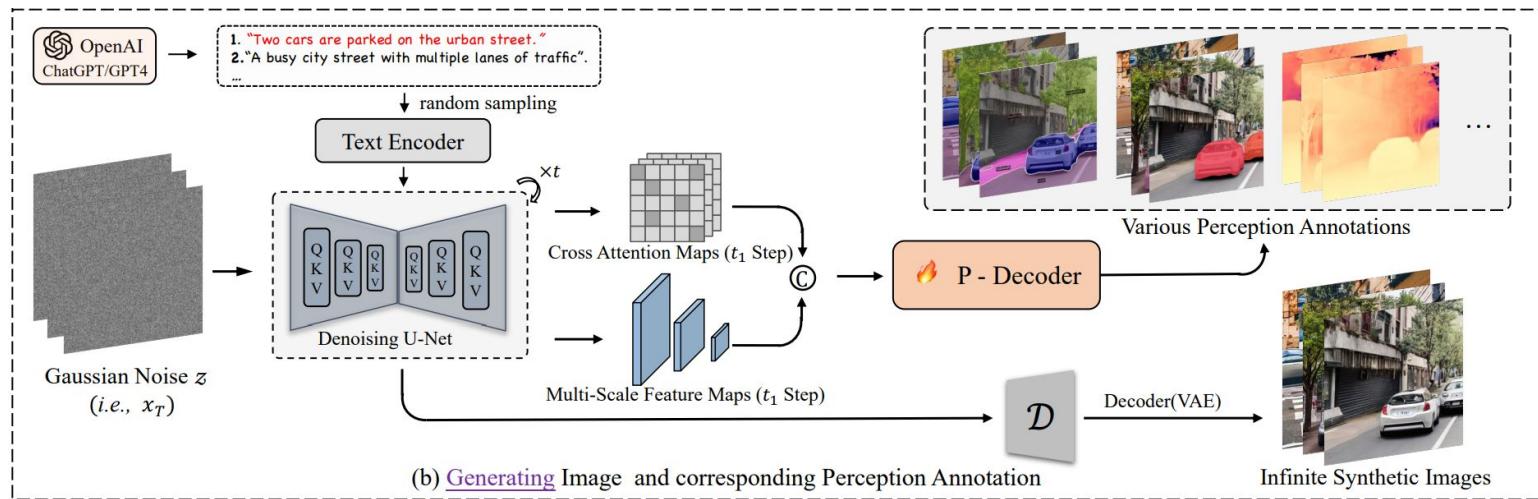
- Fine-tune a small **label decoder** on top of Diffusion Model features.



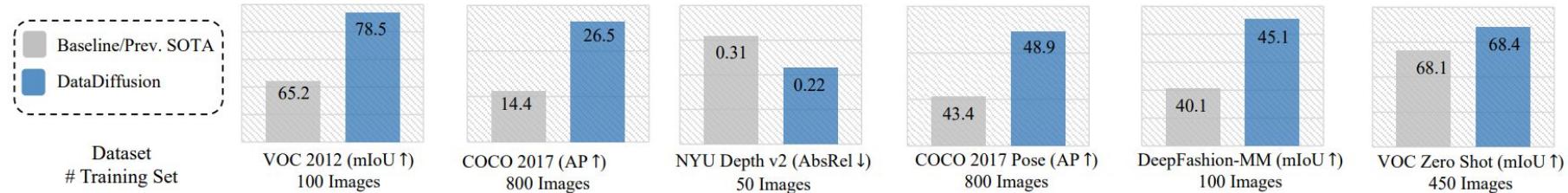
DatasetDM: Generating Images And Annotations

Key idea: the Diffusion Model has great features, just use these.

- During image generation, annotations are **very cheap to generate**.



DatasetDM: Generating Images And Annotations



- Very efficient, can generate annotations from < 1000 example images.
- Synthetic data improves vision models in data-limited settings.

DatasetDM Questions?

From Language Models To Agent Models

- A new paradigm is emerging using **Large Language Models** to automate **computer tasks** like Web Browsing, Software Dev, and more.

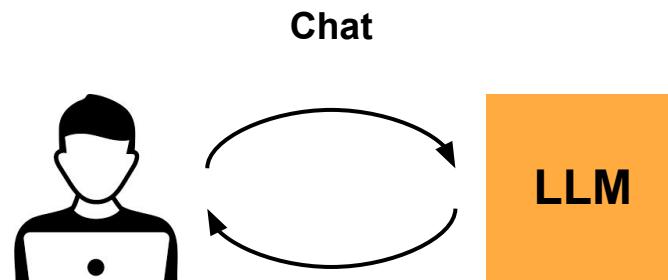
Book a trip to Seattle
on Alaska Airlines this
Saturday from SJC



From Language Models To Agent Models

- A new paradigm is emerging using **Large Language Models** to automate **computer tasks** like Web Browsing, Software Dev, and more.

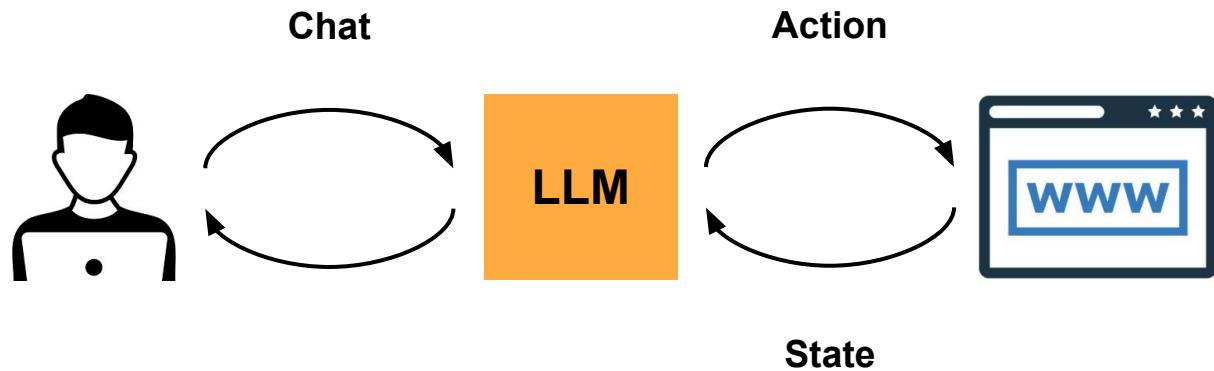
Book a trip to Seattle
on Alaska Airlines this
Saturday from SJC



From Language Models To Agent Models

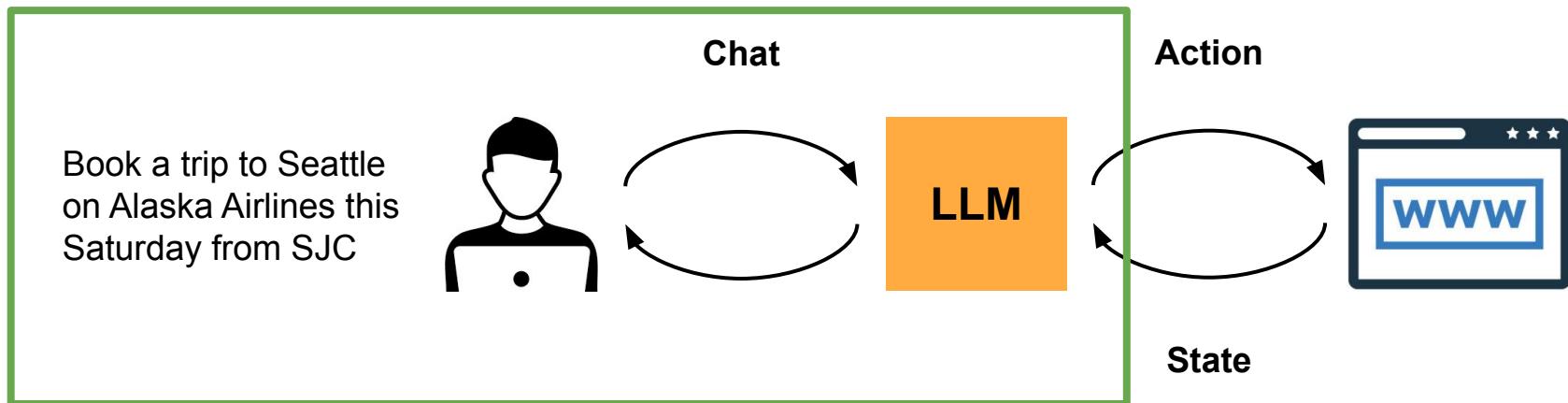
- A new paradigm is emerging using **Large Language Models** to automate **computer tasks** like Web Browsing, Software Dev, and more.

Book a trip to Seattle
on Alaska Airlines this
Saturday from SJC



From Language Models To Agent Models

- A new paradigm is emerging using **Large Language Models** to automate **computer tasks** like Web Browsing, Software Dev, and more.



Our Goal: Unlock A Critical Data Flywheel

Berkeley Library DIGITAL COLLECTIONS

Search for photos, manuscripts, newspapers and much more..

Search

Search Tips : Advanced Search

Powered by TIND

Task: Find a digitized historical map of the San Francisco Bay Area.

SCORE \$100 OFF FOR EVERY \$1000 YOU SPEND! Use Code SCORE At Checkout!

Choose from 40+ 4WP stores

Need help finding parts? Question on fitment? We can help! Give us a call 877-474-4821 or send us a message



Search...



FIND A LOCATION

SUPPORT



DIRT

STORES LOCATIONS ▾ SHOP BY CATEGORY ▾ SHOP BY BRAND ▾

JOIN US!
**SATURDAY
MARCH 8TH**

ALL LOCATIONS CELEBRATING

- FREE FOOD & DRINKS
- MUSIC, GAMES & FAMILY FUN
- DOG & FAMILY FRIENDLY

SAVE STORE-WIDE
OFF-ROAD VEHICLES ON DISPLAY
WIN FREE RAFFLE PRIZES

9AM - 3PM AT ALL NATIONWIDE LOCATIONS. 3RD - 10TH STORE-WIDE DISCOUNTS.
CALL OR VISIT YOUR NEAREST STORE FOR MORE DETAILS.

SELECT VEHICLE

My Garage ▾

Select Year

SHOP BY CATEGORY

Select Main Category

SHOP BY BRAND

Select Brand



SHOP TOP CATEGORIES

Explore promotions

Task: Find the price of a Jeep Wrangler lift kit.

Towards Internet-Scale Training For Agents

- Goal: unlock internet-scale data for training for agent models.

Towards Internet-Scale Training For Agents

- **Goal:** unlock internet-scale data for training for agent models.
- **Key Idea:** LLMs can generate and verify agent tasks on live websites.

Towards Internet-Scale Training For Agents

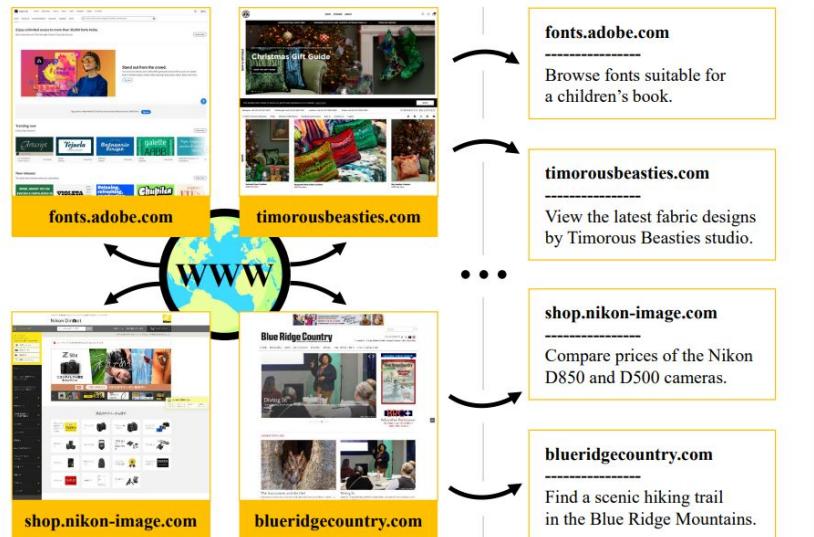
- **Goal:** unlock internet-scale data for training for agent models.
- **Key Idea:** LLMs can generate and verify agent tasks on live websites.



1,000,000 Websites

Towards Internet-Scale Training For Agents

- **Goal:** unlock internet-scale data for training for agent models.
- **Key Idea:** LLMs can generate and verify agent tasks on live websites.

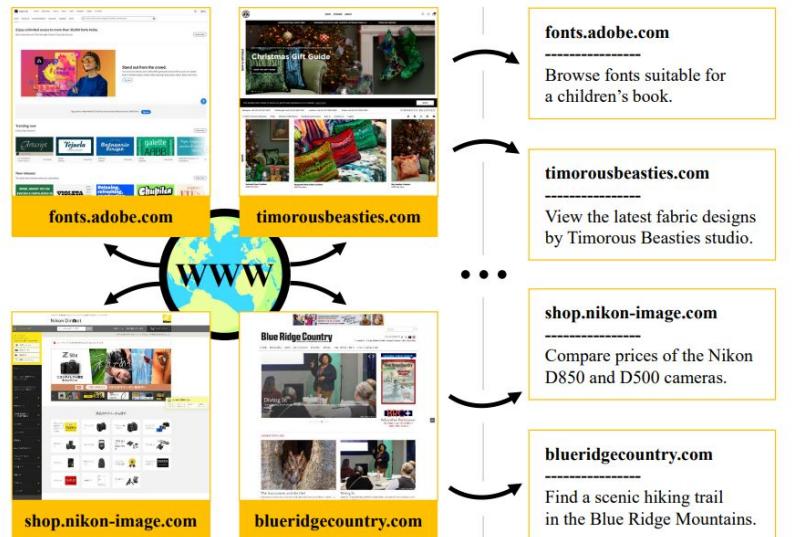


1,000,000 Websites

Stage 1: Task Generation

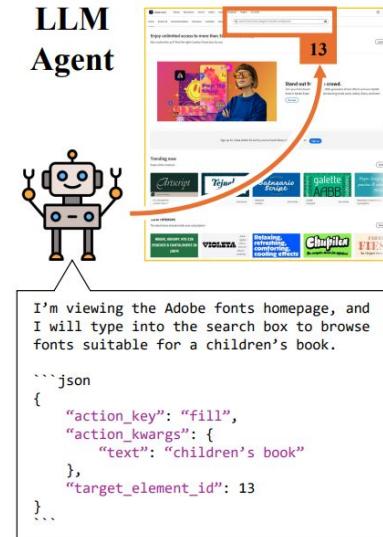
Towards Internet-Scale Training For Agents

- **Goal:** unlock internet-scale data for training for agent models.
- **Key Idea:** LLMs can generate and verify agent tasks on live websites.



1,000,000 Websites

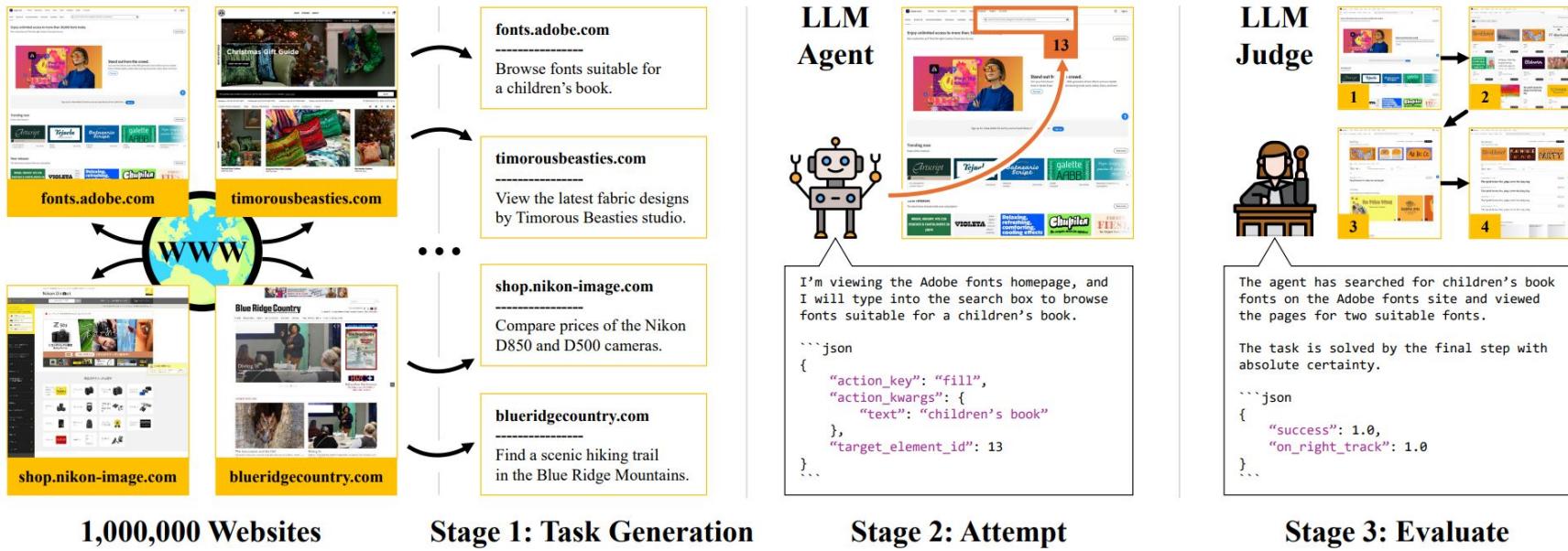
Stage 1: Task Generation



Stage 2: Attempt

Towards Internet-Scale Training For Agents

- **Goal:** unlock internet-scale data for training for agent models.
- **Key Idea:** LLMs can generate and verify agent tasks on live websites.



1,000,000 Websites

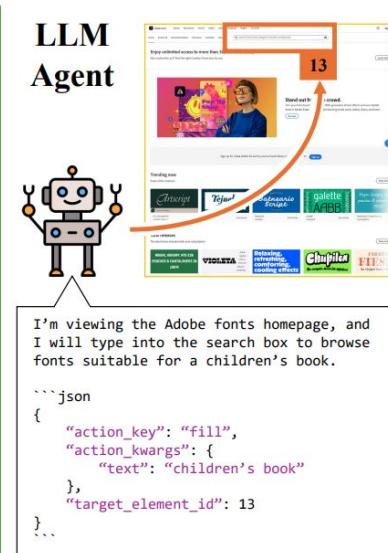
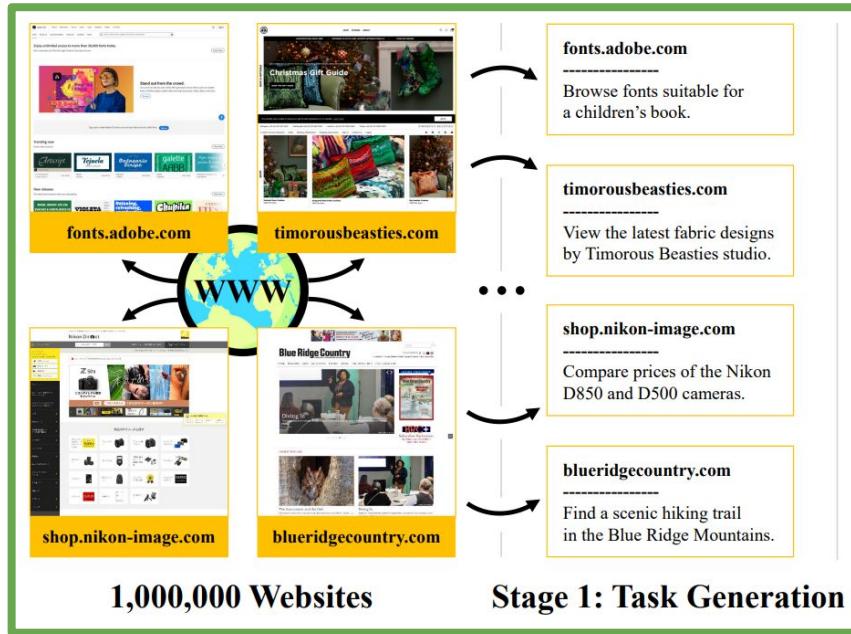
Stage 1: Task Generation

Stage 2: Attempt

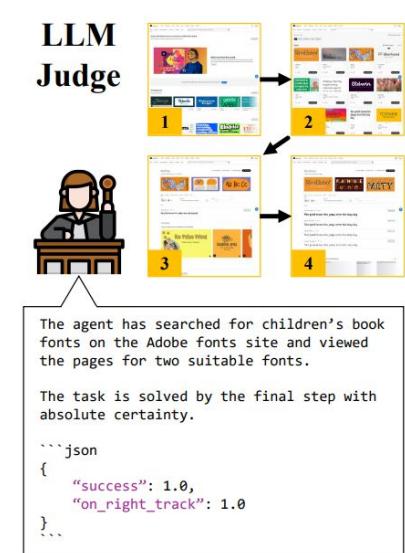
Stage 3: Evaluate

Towards Internet-Scale Training For Agents

- **Goal:** unlock internet-scale data for training for agent models.
- **Key Idea:** LLMs can generate and verify agent tasks on live websites.



Stage 2: Attempt



Stage 3: Evaluate

LLMs Can Generate **Realistic Agent Tasks**

The Problem Statement:

Given a **web domain as text** (i.e. www.uefi.org).

Propose a realistic task that **an average user could complete** in one session.

LLMs Can Generate **Realistic Agent Tasks**

The Problem Statement:

Given a **web domain as text** (i.e. www.uefi.org).

Propose a realistic task that **an average user could complete** in one session.

This is simple.

LLMs Can Generate **Realistic Agent Tasks**

The Problem Statement:

Given a **web domain as text** (i.e. www.uefi.org).

Propose a realistic task that **an average user could complete** in one session.

This is simple. Why does it work so well?

LLMs Can Generate **Realistic Agent Tasks**

The Problem Statement:

Given a **web domain as text** (i.e. www.uefi.org).

Propose a realistic task that **an average user could complete** in one session.

This is simple. Why does it work so well?

Llama was pretrained on **broad data that contains references to web domains**, and has **broad prior knowledge** about websites and their uses.

Why Is A Simple Approach So Effective?

Llama was pretrained on **broad data that contains references to web domains**, and has **broad prior knowledge** about websites and their uses.

Llama-3.1-70B-Instruct:

describe the “merseyferries.co.uk” site

Why Is A Simple Approach So Effective?

Llama was pretrained on **broad data that contains references to web domains**, and has **broad prior knowledge** about websites and their uses.

Llama-3.1-70B-Instruct:

describe the “merseyferries.co.uk” site

Mersey Ferries (merseyferries.co.uk) is the **official website for the famous Mersey Ferry service**

?

Why Is A Simple Approach So Effective?

Llama was pretrained on **broad data that contains references to web domains**, and has **broad prior knowledge** about websites and their uses.

Llama-3.1-70B-Instruct:

describe the “merseyferries.co.uk” site

Mersey Ferries (merseyferries.co.uk) is the **official website for the famous Mersey Ferry service**



Why Is A Simple Approach So Effective?

Llama was pretrained on **broad data that contains references to web domains**, and has **broad prior knowledge** about websites and their uses.

Llama-3.1-70B-Instruct: describe the “merseyferries.co.uk” site

Mersey Ferries (merseyferries.co.uk) is the **official website for the famous Mersey Ferry service** ✓ in Liverpool, England. This website serves as a one-stop platform for users to access various services and information related to the ferry operation. Here are some things a user can do on this site:



Why Is A Simple Approach So Effective?

Llama was pretrained on **broad data that contains references to web domains**, and has **broad prior knowledge** about websites and their uses.

Llama-3.1-70B-Instruct: describe the “merseyferries.co.uk” site

Mersey Ferries (merseyferries.co.uk) is the **official website for the famous Mersey Ferry service**  in Liverpool, England. This website serves as a one-stop platform for users to access various services and information related to the ferry operation. Here are some things a user can do on this site:

1. **Plan a Journey:** Users can plan their journey by checking the ferry timetable, which includes departure and arrival times for both commuter and cruise services.
2. **Book Tickets:** Users can book tickets online for both commuter and cruise services, including the popular Manchester Ship Canal Cruise and the Liverpool Bay Cruise. ...

LLMs Are Competitive With Human Annotators

We compare LLMs to (non-expert) human annotators at **detecting harmful sites**, and creating **feasible agent tasks**, for a set of 100 sites.

LLMs Are Competitive With Human Annotators

We compare LLMs to (non-expert) human annotators at **detecting harmful sites**, and creating **feasible agent tasks**, for a set of 100 sites.

Method	Acc.	Prec.	Recall
<i>Llama 3.1 70B</i>	85%	0.77	1.00
<i>GPT-4o</i>	95%	0.91	1.00
<i>Gemini 1.5 Pro</i>	97%	0.96	0.98
Human Baseline	75%	0.71	0.84

Accuracy For Detecting Harmful Sites

Method	Feasibility Rate
<i>Llama 3.1 70B</i>	75%
<i>GPT-4o</i>	85%
<i>Gemini 1.5 Pro</i>	89%
Human Baseline	54%

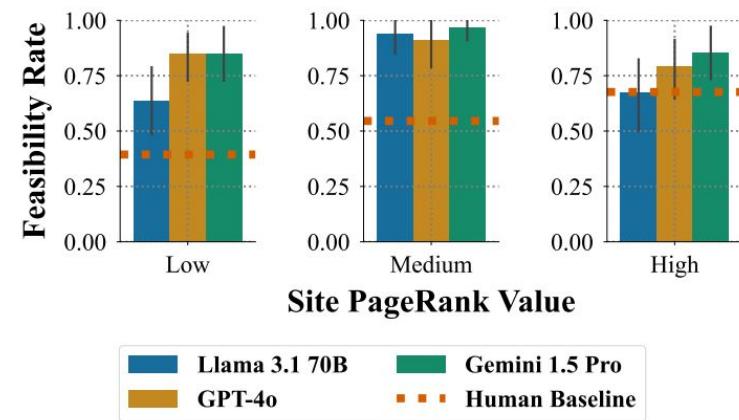
Expert Feasibility Of Proposed Tasks

LLMs Are Competitive With Human Annotators

Human annotators become **less reliable** as sites become **less popular**, and **LLMs perform consistently well** for diverse sites.

Method	Acc.	Prec.	Recall
<i>Llama 3.1 70B</i>	85%	0.77	1.00
<i>GPT-4o</i>	95%	0.91	1.00
<i>Gemini 1.5 Pro</i>	97%	0.96	0.98
Human Baseline	75%	0.71	0.84

Accuracy For Detecting Harmful Sites



Scaling To 150,000 Sites Using Llama

- Starting from the top 1M websites



1,000,000 Websites

Language Model Task Proposer

Filtered Tasks

You are helping us create tasks for a web navigation system. We will tell you the domain of a website. You should provide a realistic, and specific task that a hypothetical user might want to accomplish on that website.

Skipping Unsafe / Inappropriate Domains

To skip a domain, respond with 'N/A' instead of providing a task.

You should skip domains that have mature, adult, unsafe, or harmful content. If you are unsure whether a domain is safe, please skip it. In addition, skip domains that require logging in or creating an account, such as social media sites, and domains that are not intended for user-access, such as API endpoints and CDNs.

Here are some domains to provide tasks for:

- * www.amazon.com: 'Find the price of the 24in LG Ultragear Monitor.'
- * www.wikipedia.org: 'Look up the history of the Eiffel Tower on Wikipedia.'

Here are some domains to skip:

- * fbcdn.net: 'N/A'
- * api.github.com: 'N/A'

Tasks should not require external knowledge, not modify the state of the web, and should not require logging in or creating an account. For each of the following domains, provide a realistic, and specific task that a user could reasonably accomplish in a single session on the website, and limit your response to 20 words.

fonts.adobe.com

Browse fonts suitable for a children's book.

timorousbeasties.com

View the latest fabric designs by Timorous Beasties studio.

shop.nikon-image.com

Compare prices of the Nikon D850 and D500 cameras.

harmful-website.com



Scaling To 150,000 Sites Using Llama

- Starting from the **top 1M websites**, we generate **150k safe tasks**.



1,000,000 Websites

You are helping us create tasks for a web navigation system. We will tell you the domain of a website. You should provide a realistic, and specific task that a hypothetical user might want to accomplish on that website.

Skipping Unsafe / Inappropriate Domains

To skip a domain, respond with 'N/A' instead of providing a task.

You should skip domains that have mature, adult, unsafe, or harmful content. If you are unsure whether a domain is safe, please skip it. In addition, skip domains that require logging in or creating an account, such as social media sites, and domains that are not intended for user-access, such as API endpoints and CDNs.

Here are some domains to provide tasks for:

- * 'www.amazon.com': 'Find the price of the 24in LG Ultragear Monitor.'
- * 'www.wikipedia.org': 'Look up the history of the Eiffel Tower on Wikipedia.'

Here are some domains to skip:

- * 'fbcdn.net': 'N/A'
- * 'api.github.com': 'N/A'

Tasks should not require external knowledge, not modify the state of the web, and should not require logging in or creating an account. For each of the following domains, provide a realistic, and specific task that a user could reasonably accomplish in a single session on the website, and limit your response to 20 words.

Language Model Task Proposer



Filtered Tasks

Exploring Tasks Generated By InSTA

- Tasks are **diverse**, and many require **multiple steps of reasoning** (i.e., what makes a font **suitable for a children's book**).

Domain	Task
wordpress.org	Find a free and popular theme for a personal blog.
policies.google.com	Read Google's terms of service for using YouTube.
ec.europa.eu	Retrieve a report on the EU's climate change policy.
vimeo.com	Find a short film on environmental conservation.
fonts.adobe.com	Browse fonts suitable for a children's book.
apps.apple.com	Find the top-rated free productivity app for iPhone.

Exploring Tasks Generated By InSTA

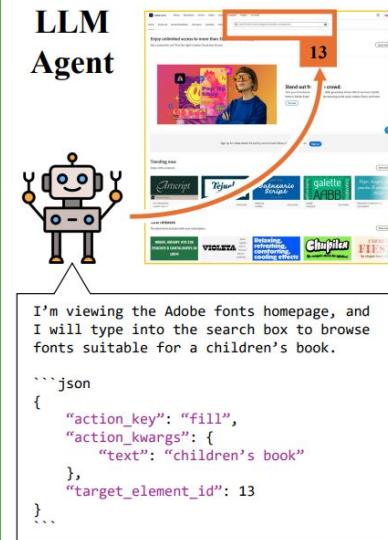
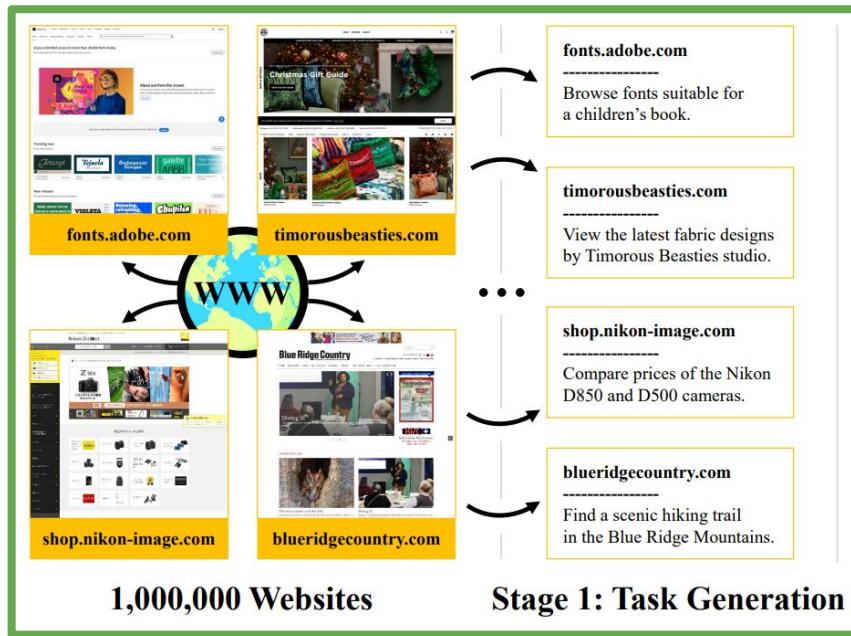
- There is an **emergent creativity** where Llama can recall facts that a site is likely to contain, such as the meaning of the Om symbol (ॐ).

Domain	Task
ancient-symbols.com	Look up the meaning of the Om symbol in ancient cultures.
petsforhomes.com.au	Find a list of available dogs for adoption in New South Wales.
timorousbeasties.com	View the latest fabric designs by the Timorous Beasties studio.
shop.nikon-image.com	Compare prices of the Nikon D850 and D500 cameras.
blueridgecountry.com	Find a scenic hiking trail in the Blue Ridge Mountains.
awg-fittings.com	Find the dimensions of a 1/2\" NPT fitting.

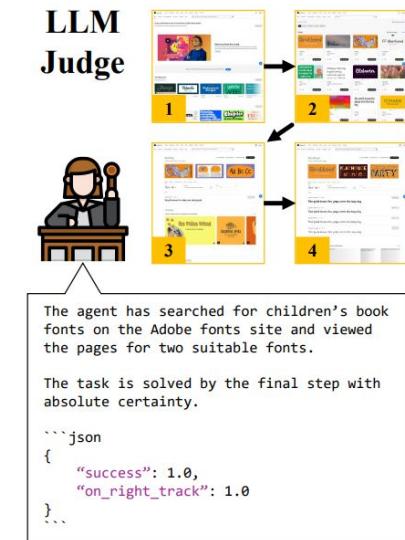
Task Generation Questions

Revisiting The InSTA Pipeline

- **Key Idea:** LLMs can **generate and verify** agent tasks on live websites.



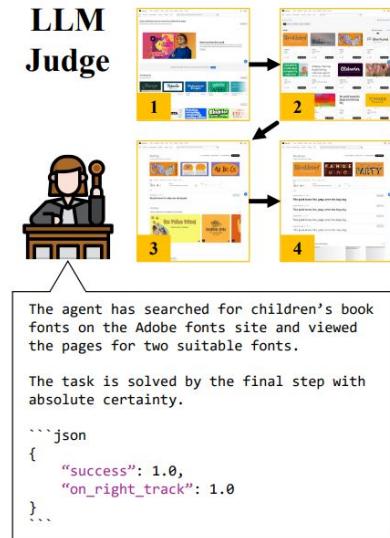
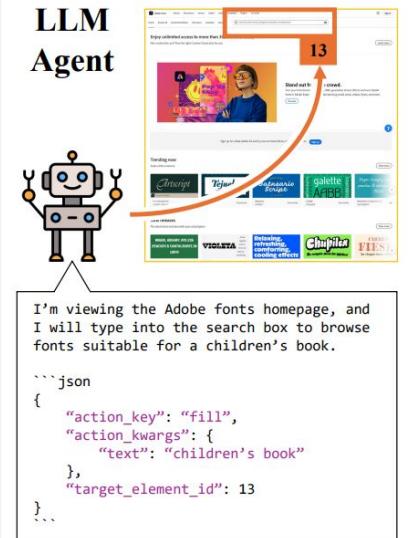
Stage 2: Attempt



Stage 3: Evaluate

Revisiting The InSTA Pipeline

- **Key Idea:** LLMs can **generate and verify** agent tasks on live websites.
- We now have **broad and diverse** agent tasks.



1,000,000 Websites

Stage 1: Task Generation

Stage 2: Attempt

Stage 3: Evaluate

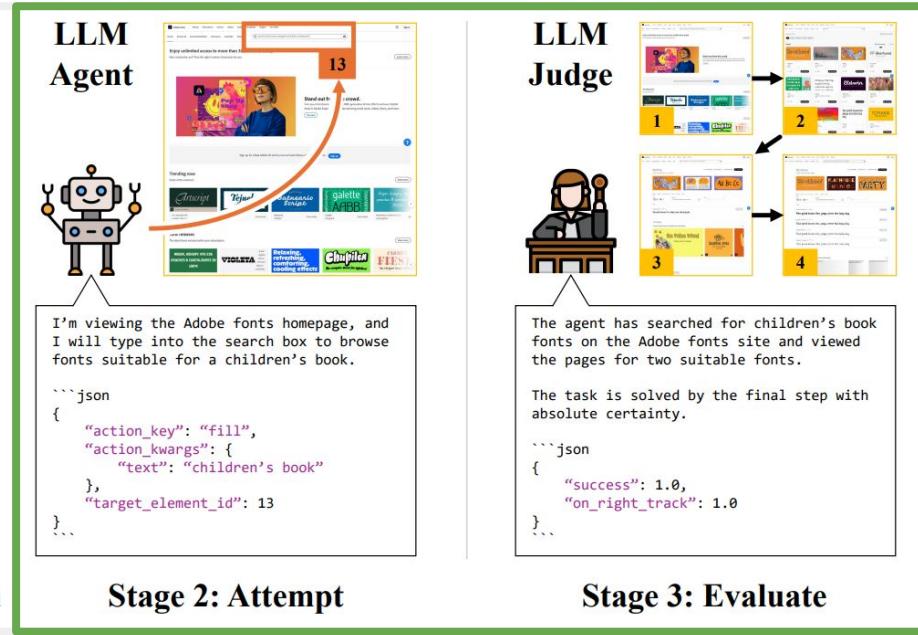
Revisiting The InSTA Pipeline

- **Key Idea:** LLMs can **generate and verify** agent tasks on live websites.
- Now we can **evaluate agents** on these tasks.



1,000,000 Websites

Stage 1: Task Generation



How Do We Evaluate Open-Ended Tasks?

How Do We Evaluate Open-Ended Tasks?

A simple strategy works for evaluation:

$$\mathbf{r}_T = f^{\text{text} \rightarrow \text{val}}(\text{LLM}([\mathbf{y}_{\text{sys}}, \mathbf{c}, \mathbf{a}_1, \dots, \mathbf{a}_T, \text{Enc}(\mathbf{s}_T)]))$$

How Do We Evaluate Open-Ended Tasks?

A simple strategy works for evaluation:

$$\mathbf{r}_T = f^{\text{text} \rightarrow \text{val}}(\text{LLM}([\mathbf{y}_{\text{sys}}, \mathbf{c}, \mathbf{a}_1, \dots, \mathbf{a}_T, \text{Enc}(\mathbf{s}_T)]))$$

- Given a task \mathbf{c} , and a trajectory $\mathbf{a}_1, \dots, \mathbf{a}_T, \text{Enc}(\mathbf{s}_T)$, we can directly estimate the probability of success via next-token prediction.

How Do We Evaluate Open-Ended Tasks?

A simple strategy works for evaluation:

$$\mathbf{r}_T = f^{\text{text} \rightarrow \text{val}}(\text{LLM}([\mathbf{y}_{\text{sys}}, \mathbf{c}, \mathbf{a}_1, \dots, \mathbf{a}_T, \text{Enc}(\mathbf{s}_T)]))$$

- Given a task \mathbf{c} , and a trajectory $\mathbf{a}_1, \dots, \mathbf{a}_T, \text{Enc}(\mathbf{s}_T)$, we can **directly estimate the probability of success** via **next-token prediction**.
- Recent work shows this **works well for math tasks** [3].

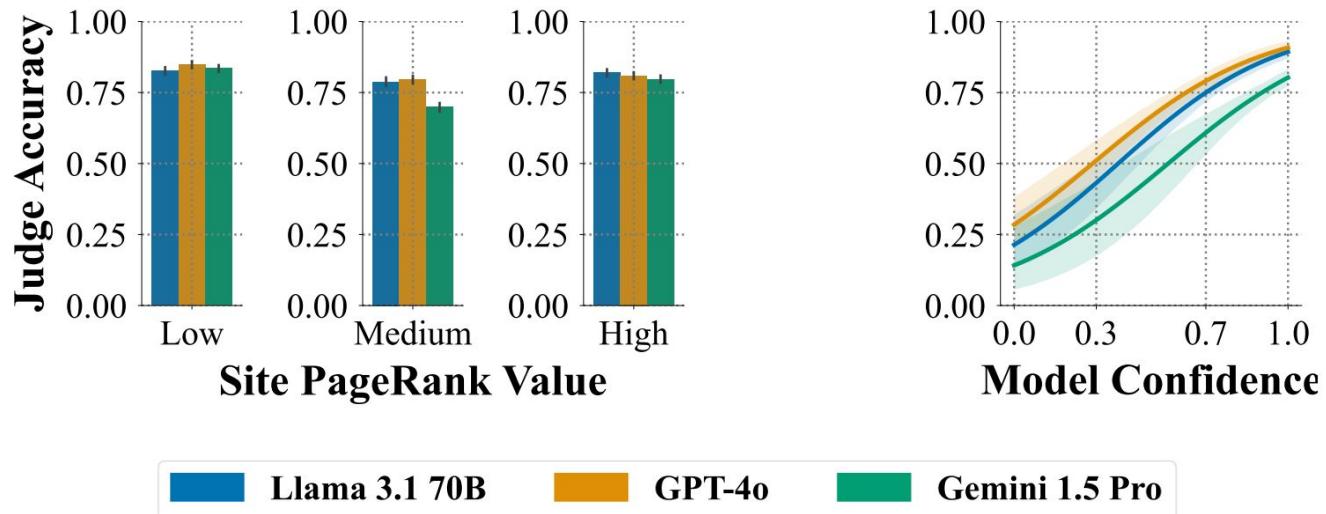
How Do We Evaluate Open-Ended Tasks?

A simple strategy works for evaluation:

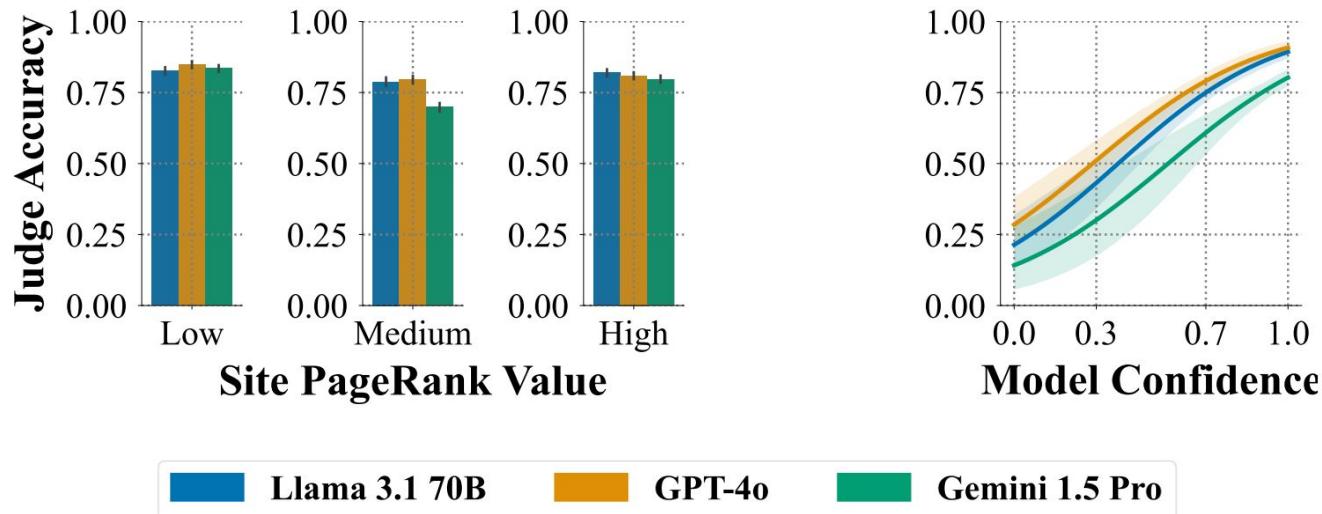
$$\mathbf{r}_T = f^{\text{text} \rightarrow \text{val}}(\text{LLM}([\mathbf{y}_{\text{sys}}, \mathbf{c}, \mathbf{a}_1, \dots, \mathbf{a}_T, \text{Enc}(\mathbf{s}_T)]))$$

- Given a task \mathbf{c} , and a trajectory $\mathbf{a}_1, \dots, \mathbf{a}_T, \text{Enc}(\mathbf{s}_T)$, we can **directly estimate the probability of success** via **next-token prediction**.
- How robust are LLMs for **evaluating agent tasks**?

Language Models Are Robust Evaluators

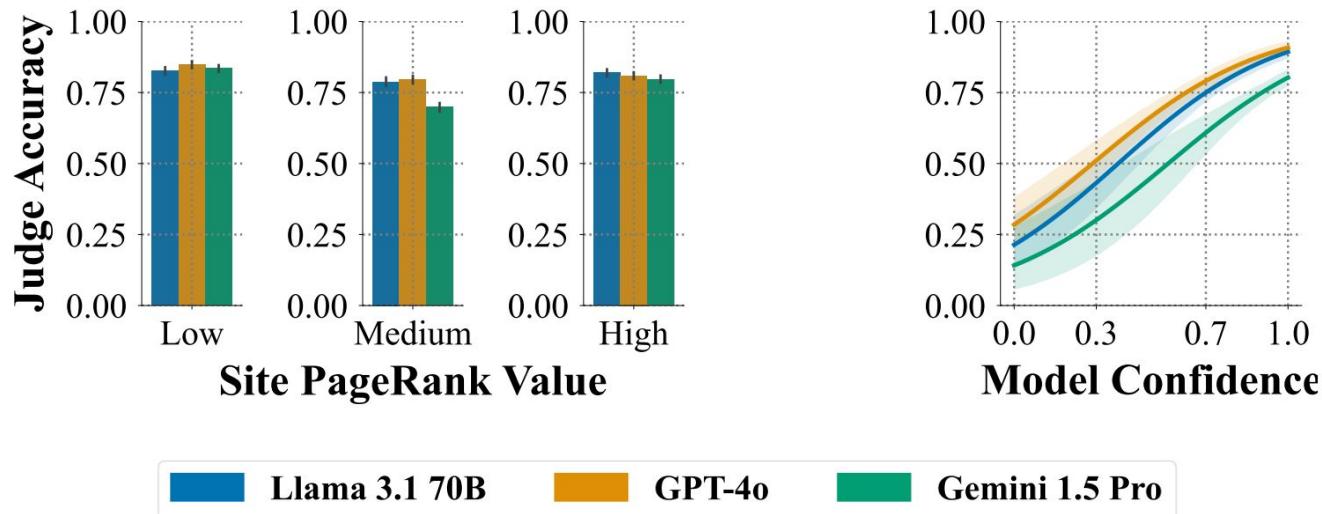


Language Models Are Robust Evaluators



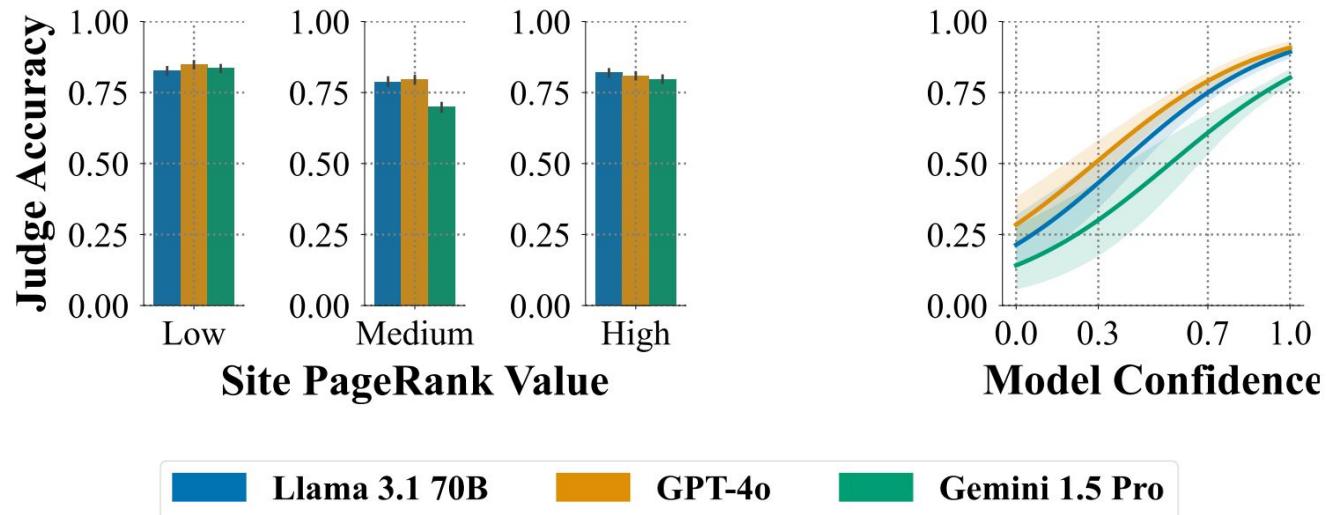
- LLMs show **high accuracy** for **detecting successful trajectories**.

Language Models Are Robust Evaluators



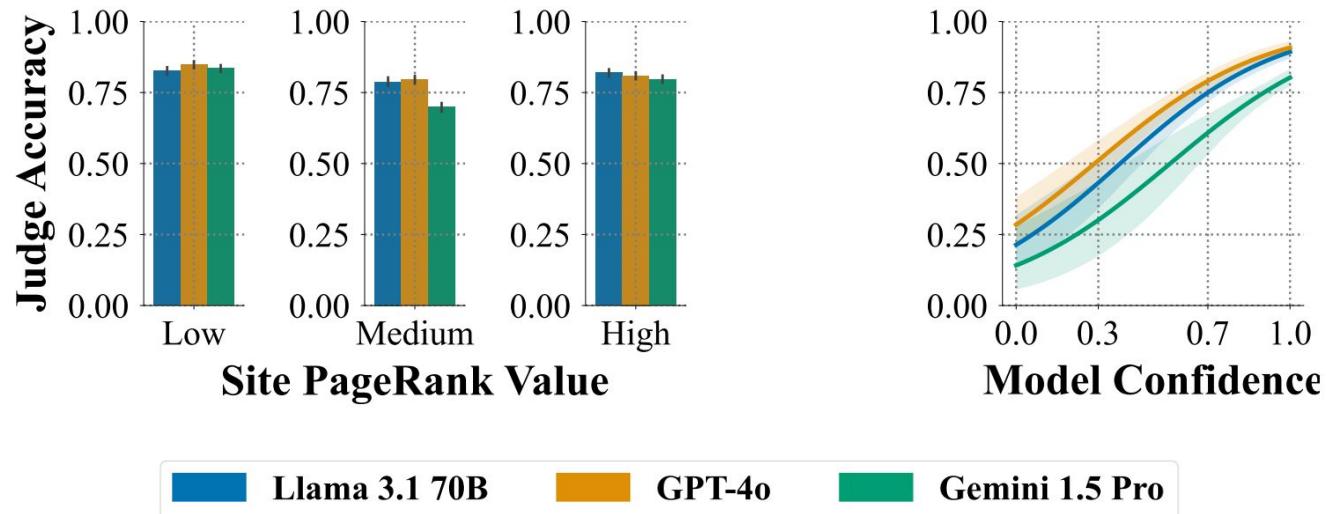
- LLMs show **high accuracy** for **detecting successful trajectories**.
- **GPT-4o** □ 82.6% , Llama 3.1 70B □ 81.7% , Gemini 1.5 Pro □ 78.0%

Language Models Are Robust Evaluators



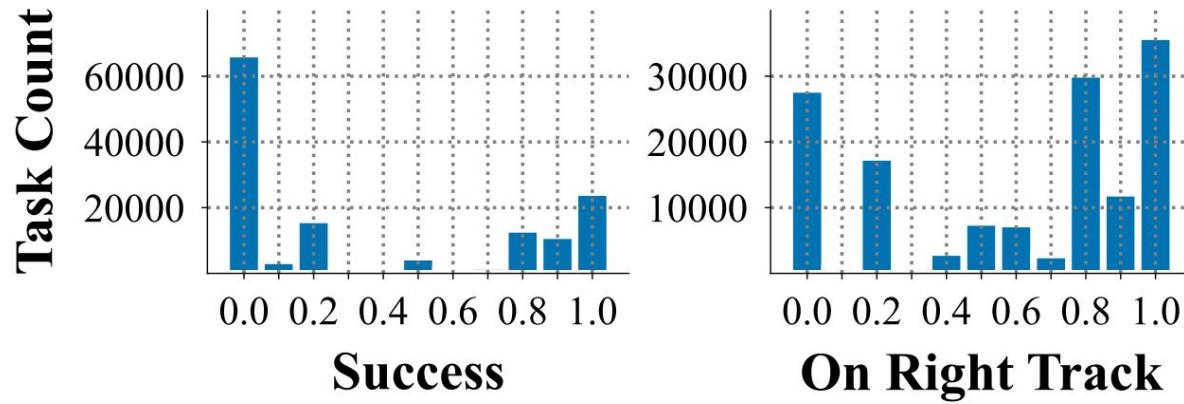
- Accuracy improves as LLMs become more confident.

Language Models Are Robust Evaluators

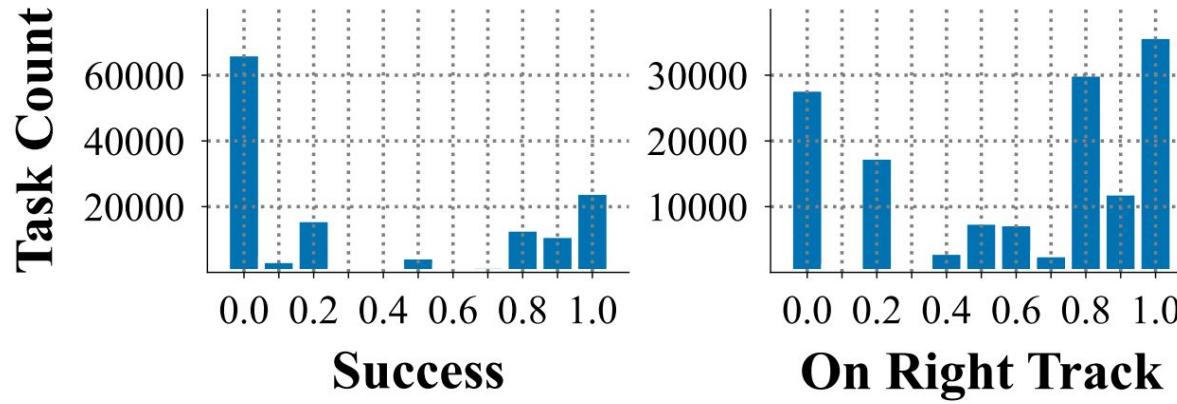


- **Accuracy improves** as LLMs become **more confident**.
- For sites with **confidence = 1.0**, **Llama 3.1 70B** **93.1% Accuracy**.

Scaling To 150,000 Agents Using Llama

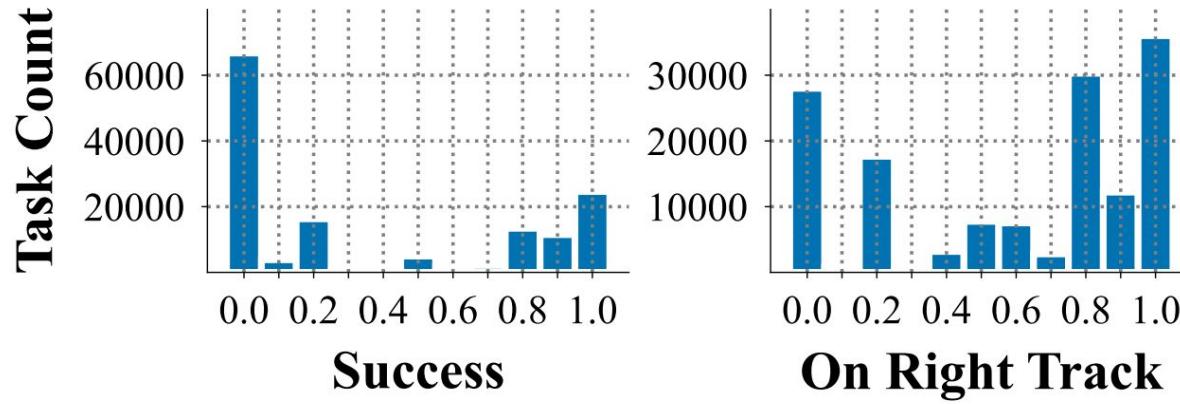


Scaling To 150,000 Agents Using Llama



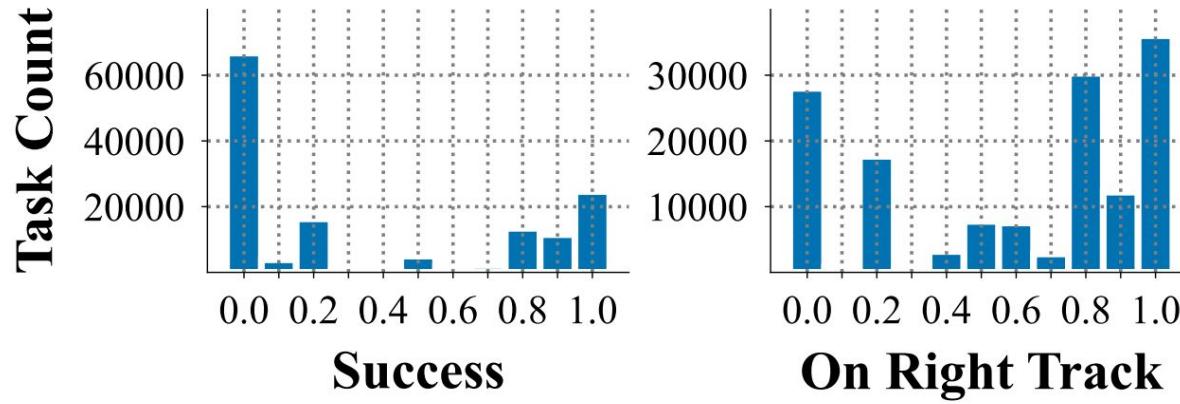
- **16.7% of tasks** are estimated to be **successful with confidence = 1.0**.

Scaling To 150,000 Agents Using Llama



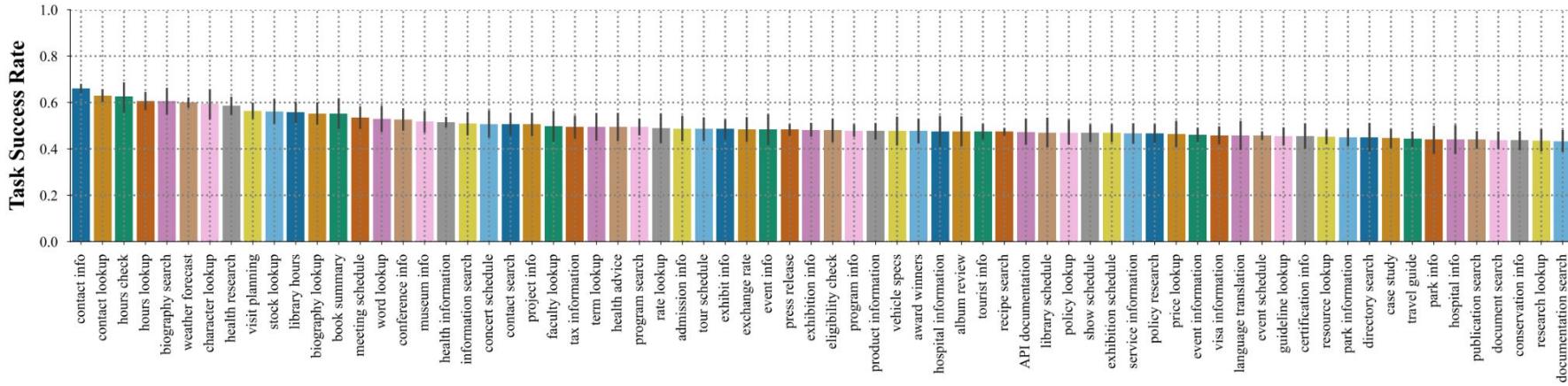
- **16.7% of tasks** are estimated to be **successful with confidence = 1.0**.
- Certain agents that failed are **on the right track**.

Scaling To 150,000 Agents Using Llama

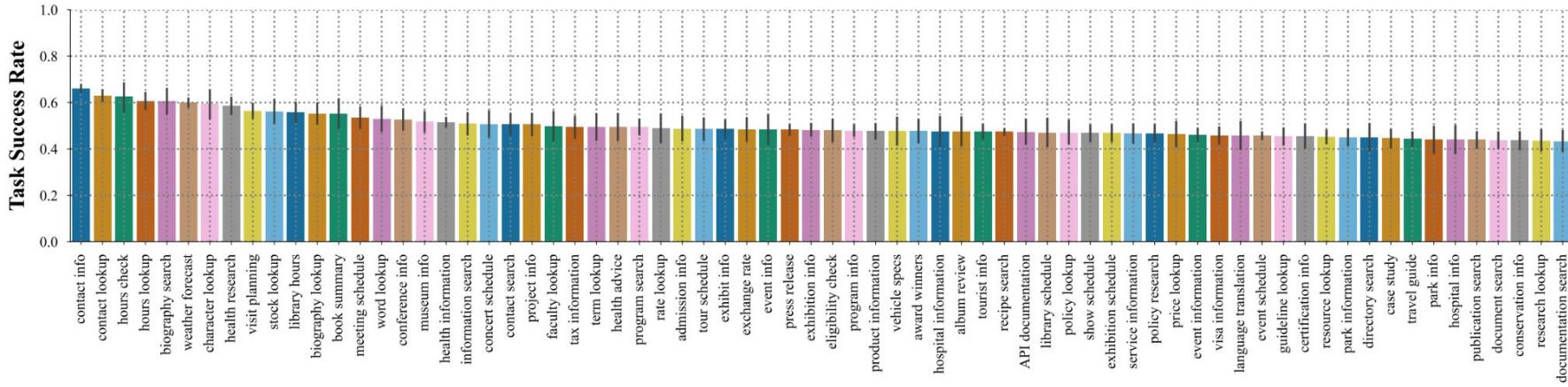


- **16.7% of tasks** are estimated to be **successful with confidence = 1.0**.
- Certain agents that failed are **on the right track**.
- Suggests **more agents can succeed** given more compute.

Most Successful Task Categories

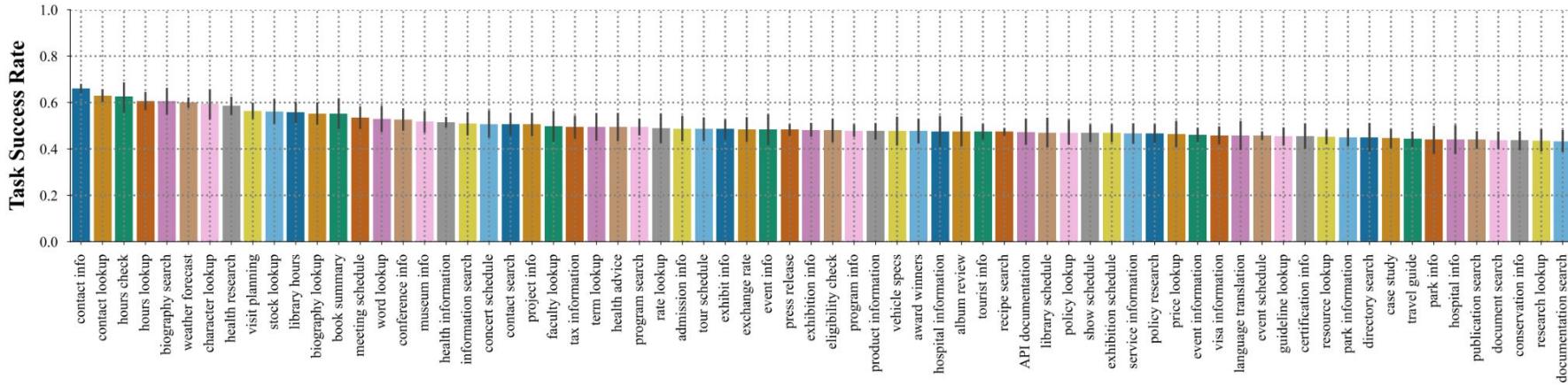


Most Successful Task Categories



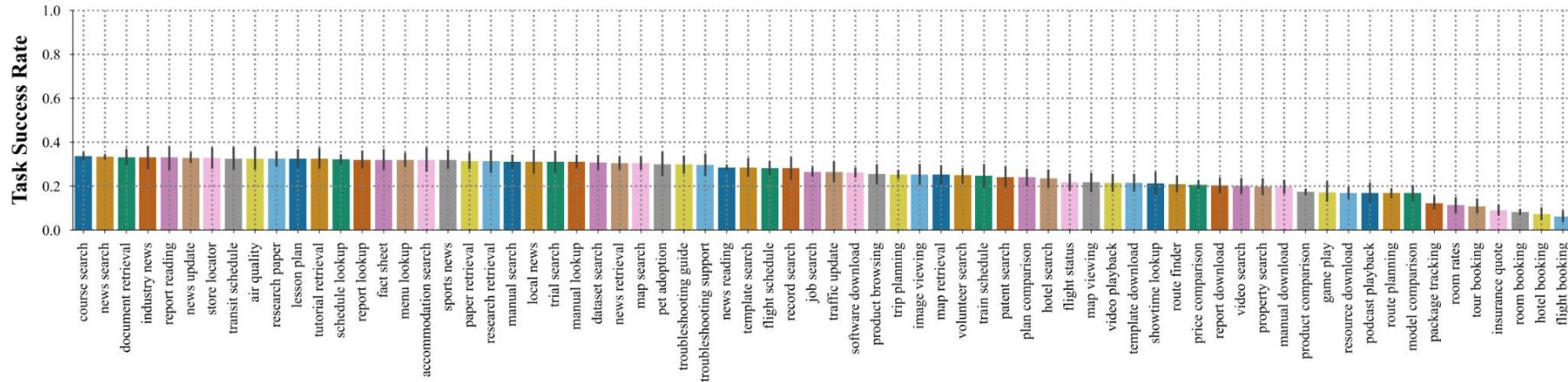
- Top 22 categories are solved with **> 50% success rate** with Llama 3.1 70B.

Most Successful Task Categories

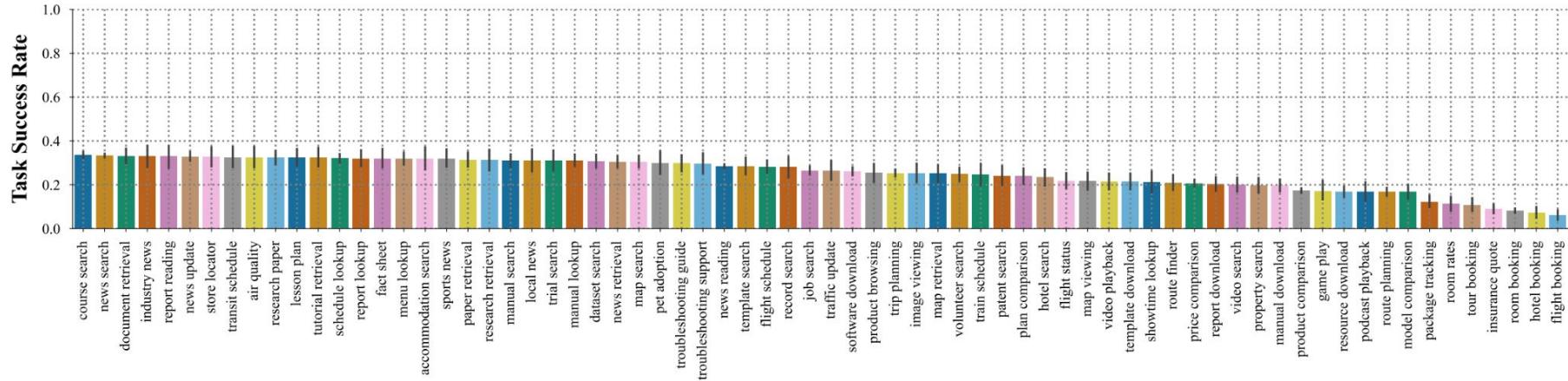


- Top 22 categories are solved with **> 50% success rate** with Llama 3.1 70B.
- Llama is most effective at **research-based tasks**.

Least Successful Task Categories

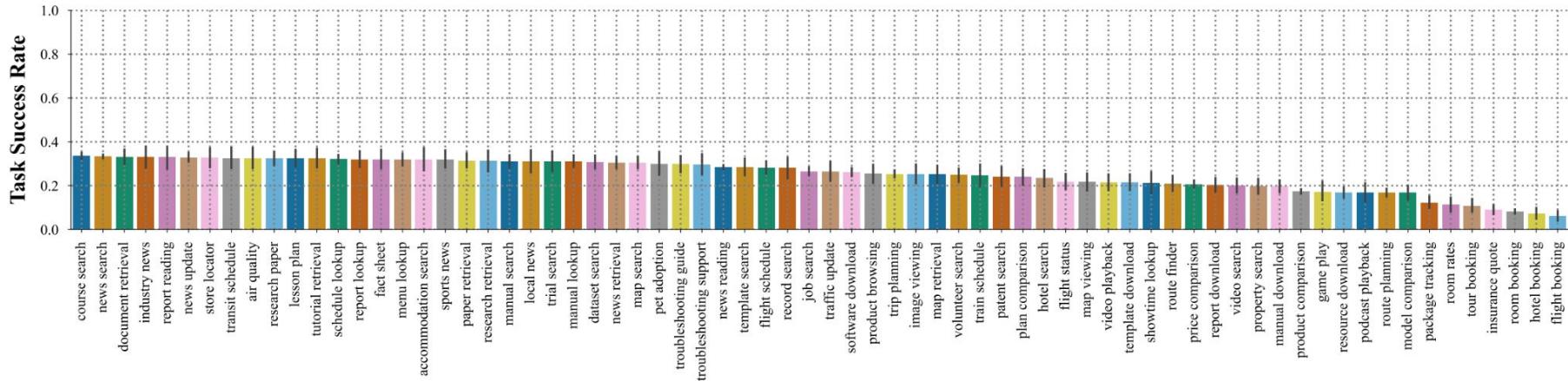


Least Successful Task Categories



- Agents often struggle to **remember and reason** about previous interactions.

Least Successful Task Categories

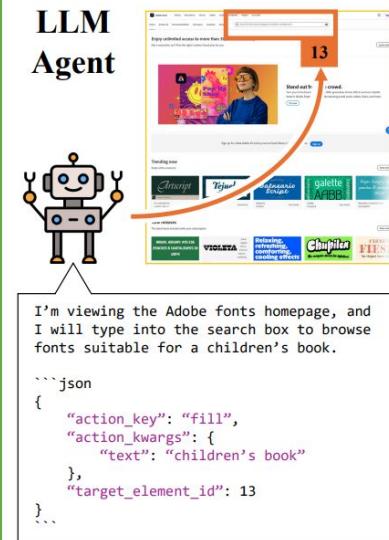
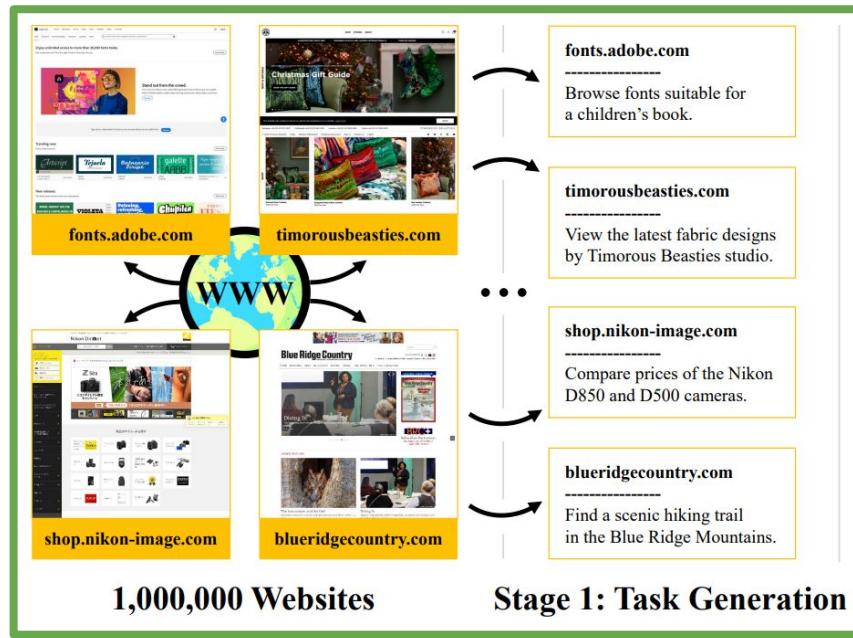


- Agents often struggle to **remember and reason** about previous interactions.
- Providing **tools like a note-taking app** could unlock these abilities.

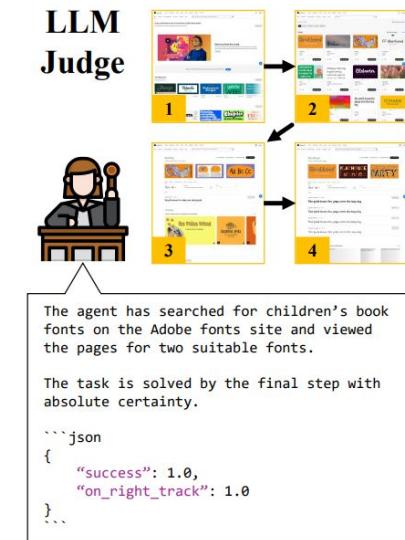
Agent & Judge Questions

Wrapping Up The InSTA Pipeline

- We've covered **generation** of agent tasks.



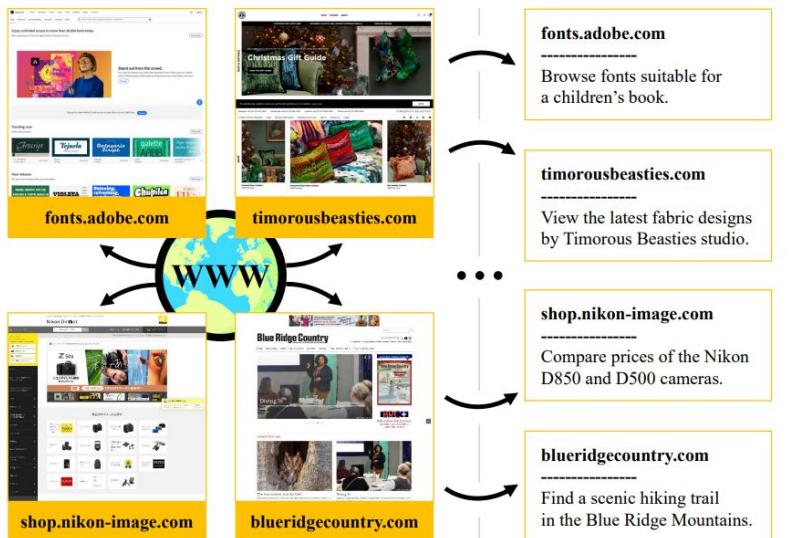
Stage 2: Attempt



Stage 3: Evaluate

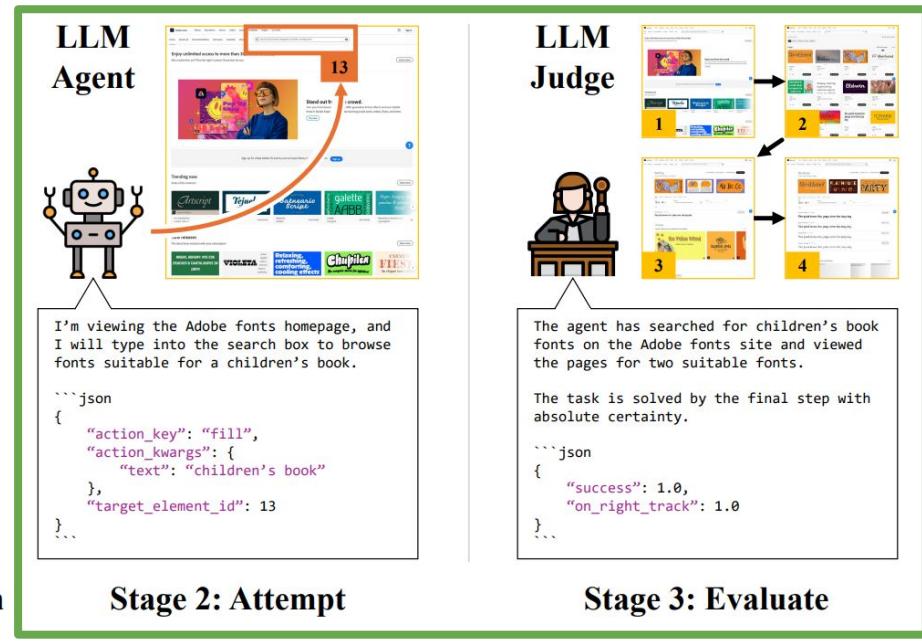
Wrapping Up The InSTA Pipeline

- We've covered **generation** and **verification** of agent tasks.



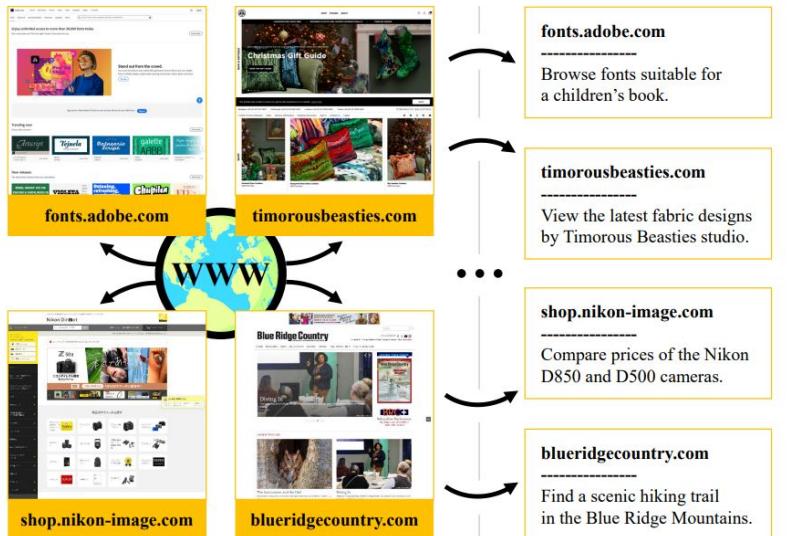
1,000,000 Websites

Stage 1: Task Generation

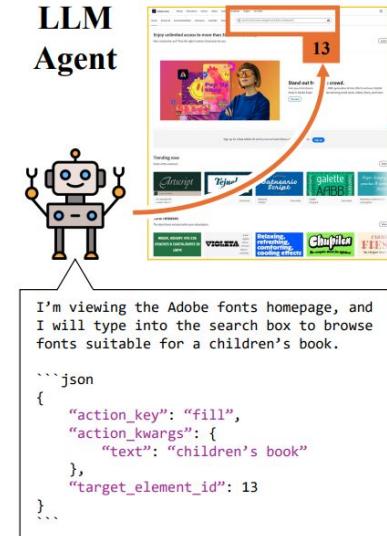


Wrapping Up The InSTA Pipeline

- We've covered **generation and verification** of agent tasks.
- How useful is this data for **training agents?**



Stage 1: Task Generation

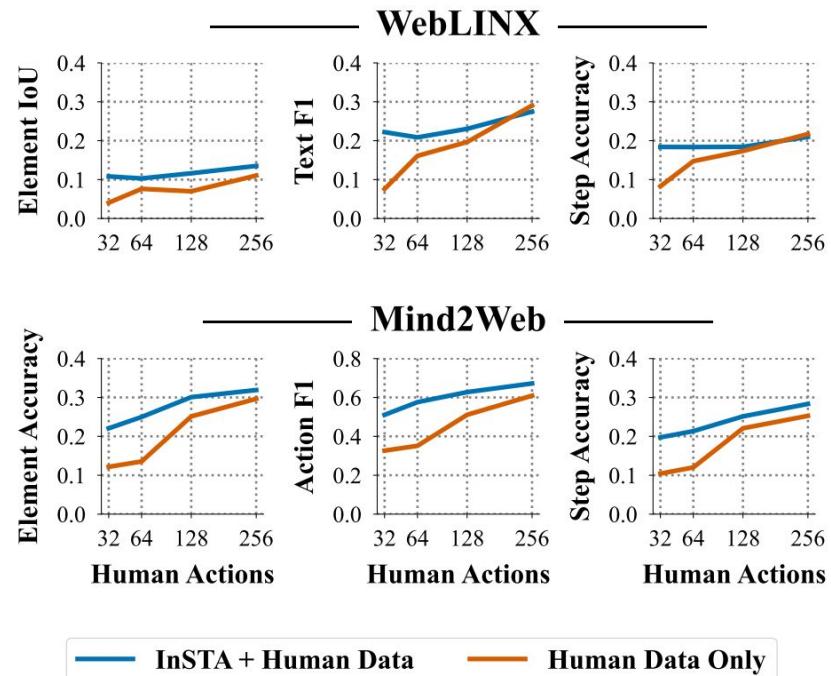


Stage 2: Attempt



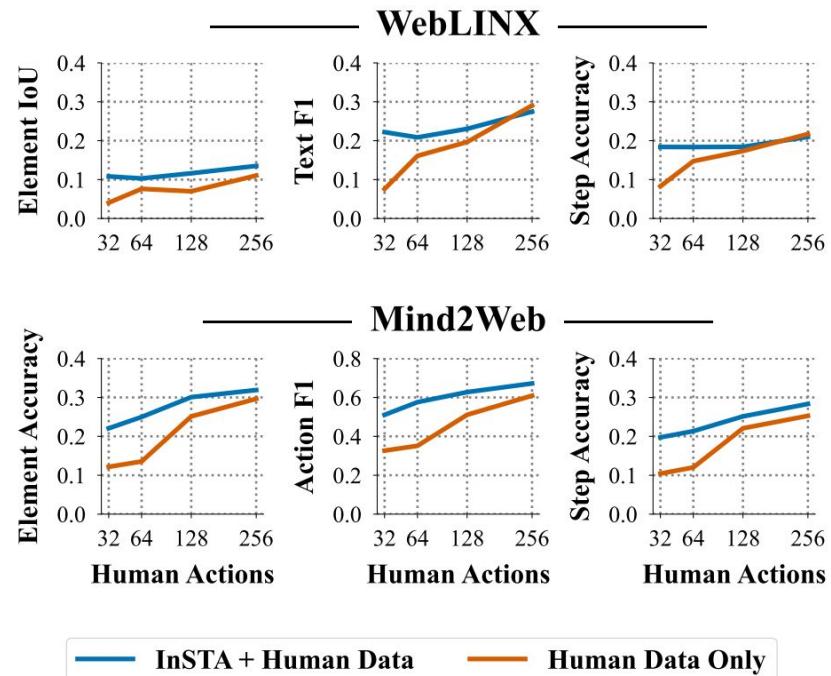
Stage 3: Evaluate

InSTA Improves Data Efficiency



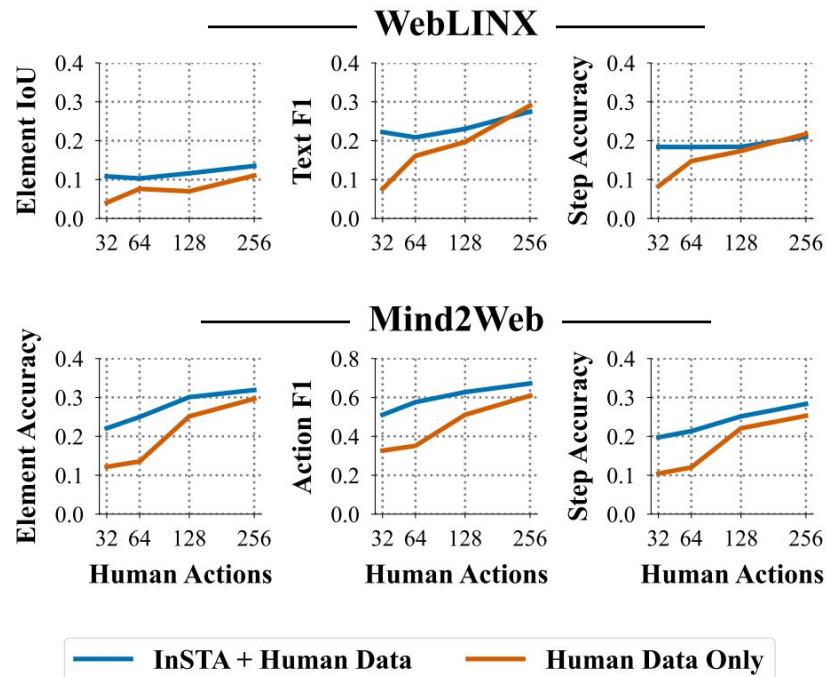
InSTA Improves Data Efficiency

- In data-limited settings derived from Mind2Web and WebLINX, SFT on a mix of our data and human data **improves efficiency**.

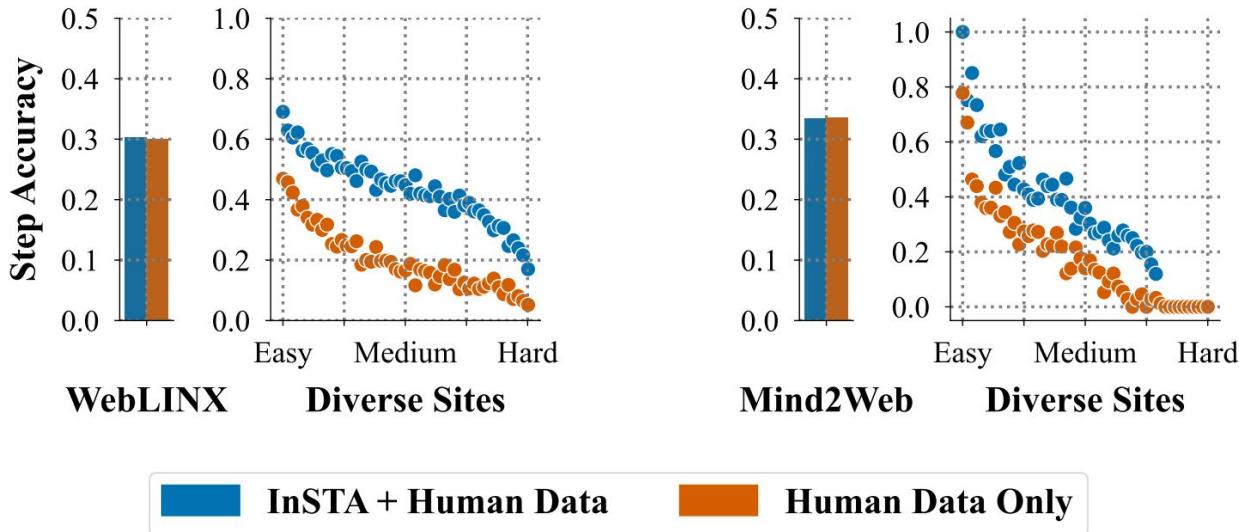


InSTA Improves Data Efficiency

- In data-limited settings derived from Mind2Web and WebLINX, SFT on a mix of our data and human data **improves efficiency**.
- Best gains for 32 human actions, performance **begins to converge** as human data size increases.

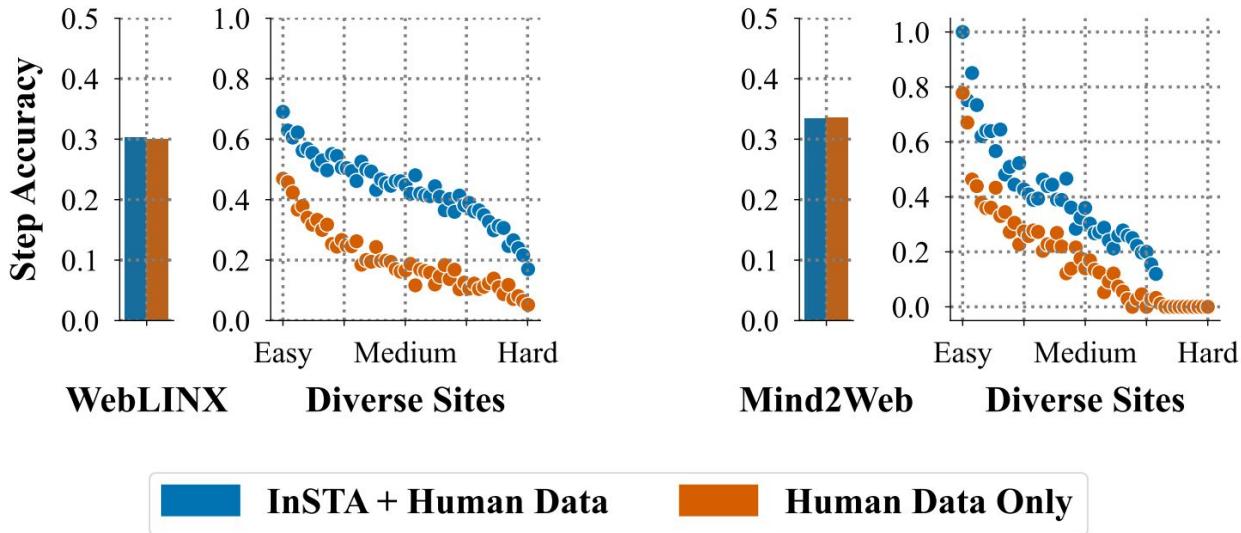


InSTA Improves Generalization



- These benchmarks test agents on a **limited set of popular sites**.

InSTA Improves Generalization



- These benchmarks test agents on a **limited set of popular sites**.
- Agents trained on **all human data** from these benchmarks **struggle to generalize** to a more diverse test set of 500 real sites.

Wrapping Up The Talk

- We developed InSTA, a pipeline **driven by language models** to unlock **internet-scale data** for training agent models.
- Scaled agents to 150k live web tasks.
- InSTA can improve **data-efficiency**, and **generalization**.

Wrapping Up The Talk

- We developed InSTA, a pipeline **driven by language models** to unlock **internet-scale data** for training agent models.
- Scaled agents to 150k live web tasks.
- InSTA can improve **data-efficiency**, and **generalization**.

What Comes Next?

There are **385M unique domains** in the Common Crawl PageRank [4], so a significant amount of data is left to **scale further**.

Wrapping Up The Talk

- We developed InSTA, a pipeline **driven by language models** to unlock **internet-scale data** for training agent models.
- Scaled agents to 150k live web tasks.
- InSTA can improve **data-efficiency**, and **generalization**.

What Comes Next?

There are **385M unique domains** in the Common Crawl PageRank [4], so a significant amount of data is left to **scale further**.

I'm excited about Reinforcement Learning with InSTA.

Final Questions?