# MULTIMOIDAL DEEPRESEARCHER: GENERATING TEXT-CHART INTERLEVED REPORTS FROM SCRATCH WITH AGENTIC FRAMEFWORK
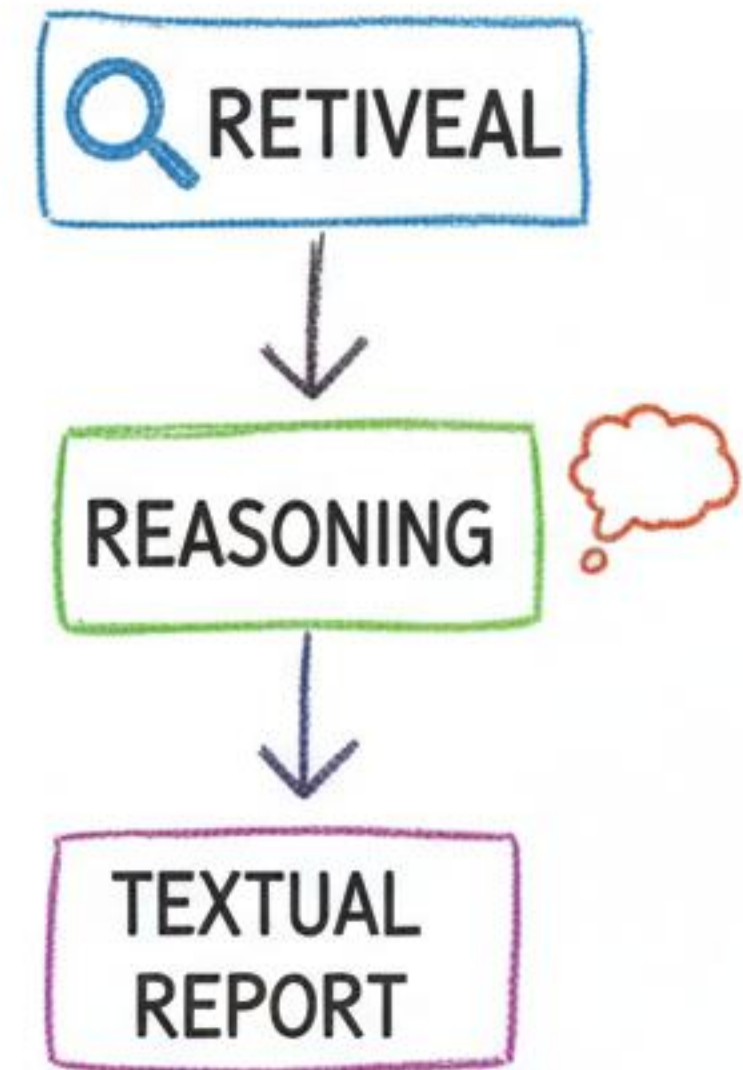
PRESENTER: SHIN-CHAN NOHARA
AFFILIATION: FUTABA KINDGERTEN (FILLER TEXT)

# RESEARCH BACKGROUND: FROM TEXTUAL REPORTS TO MULTIDODAL DEEP INQUIRY

- LARGE LANGUAGE MODELS (LLMS) EXCEL IN QA, CODE GENERATION, AND MATH.

- RETRIEVAL-AUGUNTED RESEARCH FRAMEWORKS ENABLE LOWUSEE EXTERNAL KNOWLEDGE FOR REPORTS.

- CURRENT DEEP INQUIRY FRAMEWORKS (ACADEMIC & INDUSTRIAL) MAINLY TO MAINLY PRODUCE TEXT-ONLY REPORTS, IGNORING VISUALSZATION'S ROLE.

- TEXT-HEAVY, CHART-LESS REPORTS HINDER PATTERN PACCOVERY, INFORMATION ABSORPTION, AND AUDIENCE ENGAGEMENT

**TRADITIONAL DEEP INQUIRY**



**VISUALS ARE KEY FOR BETTER UNDERSTANDING & ENGAGEMENT!**
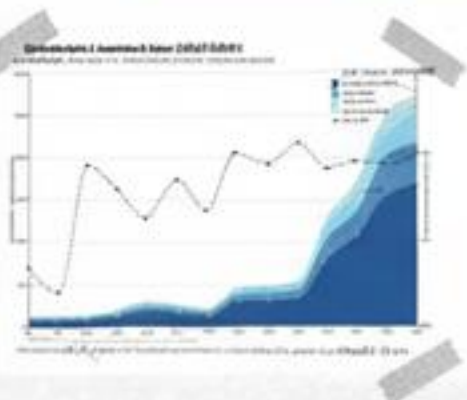
# RESEARCH QUESTIONS AND CONTRIBUTIONS OVERVIEW

## CORE PROBLEMS

How to automaticaly generate multi-modal Gesearch with intrsven text & charts from, beyond single text?

**Challenge 1:** Lack to unifed, structured chart description forms, in-concxt learning with with example reports diffult.

**Challenge 2:** How to plan report structure & visualation based on multi-round retrieval & reasning, reasning, maintaning overall consistenc?

**Challenge 3:** How automaticaly achieeve high-fidelity chart the task mis shases: four phases & iterative optimization, optimization, approcing human expert level?

## MAIN CONTRIBUTIONS

Proposing a new task "Multi-modal Report Generation" with corrospnding dataset & metrics (MultomodalReportBench).

Introducing Formal Description of Visualization (FDV), a a structured text represetation for arbiteray visusalization design.

Desingng the Multimodal DeepResarcher agent framework framewzohs, Example Textualizan, Planning & Planning, and Multi-modal Genodal Generation

Achieved 82% overall win rate against the modifed DataNarrtive baseline und equivalent model (Claude 3.7 Sonnet)

**Overall Framework: Multinnodal DeepRelsercher**

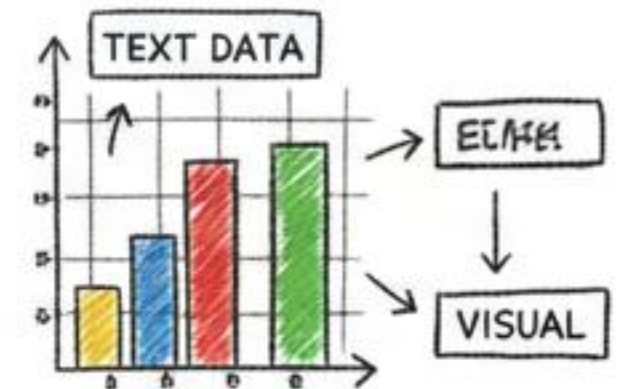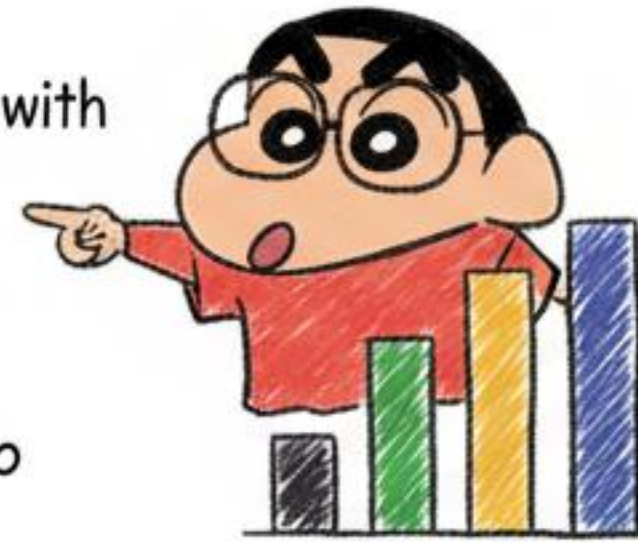# RELATED WORK: DEEP RESEARCH & VISUALIZATION GENERATION

## DEEP RESEARCH

- Combining Retireval & Reasning to push LLMs beyond parametic knowledge (e.g. OpenResuarcher, Search-o1).

- Methods use specific prompts & workflows for multi-stage reasning & retireval.

- Some explore RL end-only output, lacking multtidodal reports & chart integration.

Retrieval-Reassing Loop

## LLM FOR DATA VISULIZATION

- Focus on single chart quality: multi-stage pipelines, iterraive dbbuging with with visual feedback, , CoT-guided query refoflulzation.

- Multumodal prompting, interactive ops, multi-lingrfaces, converstional vis generation

- Eval methods mostly for single/limited chart types (bar, line) hard to support complex reports

- THIS WORK DIFFERS: First to focus on 'Text-Chart Interwoven Reports' holistic holistic gen & eval

TEXT DATA

VISUAL

Single Chart Example

# MULTIMODAL DEEPRESEARCHER: A FOUR-PHASE PROCESS



**Multimodal DeepResarcher Pipeline**

**1. RESEARCHING**

Input: Topic (t), Multinoddal
Process: Multi-turn Retiieval &
Output: Struclar: Structured
"Learnings" (L)

**2. EXEPMITUIZATION**

Input: Examples (R)
Process: Convert R to Text (uR).
Output: Steport R to Text
Eeamples)

**3. PLANNING**

Input: L, t, uR)
Process: Outline & Style Guide
Generation Visual Guide (G)
Output: Examples (uR)

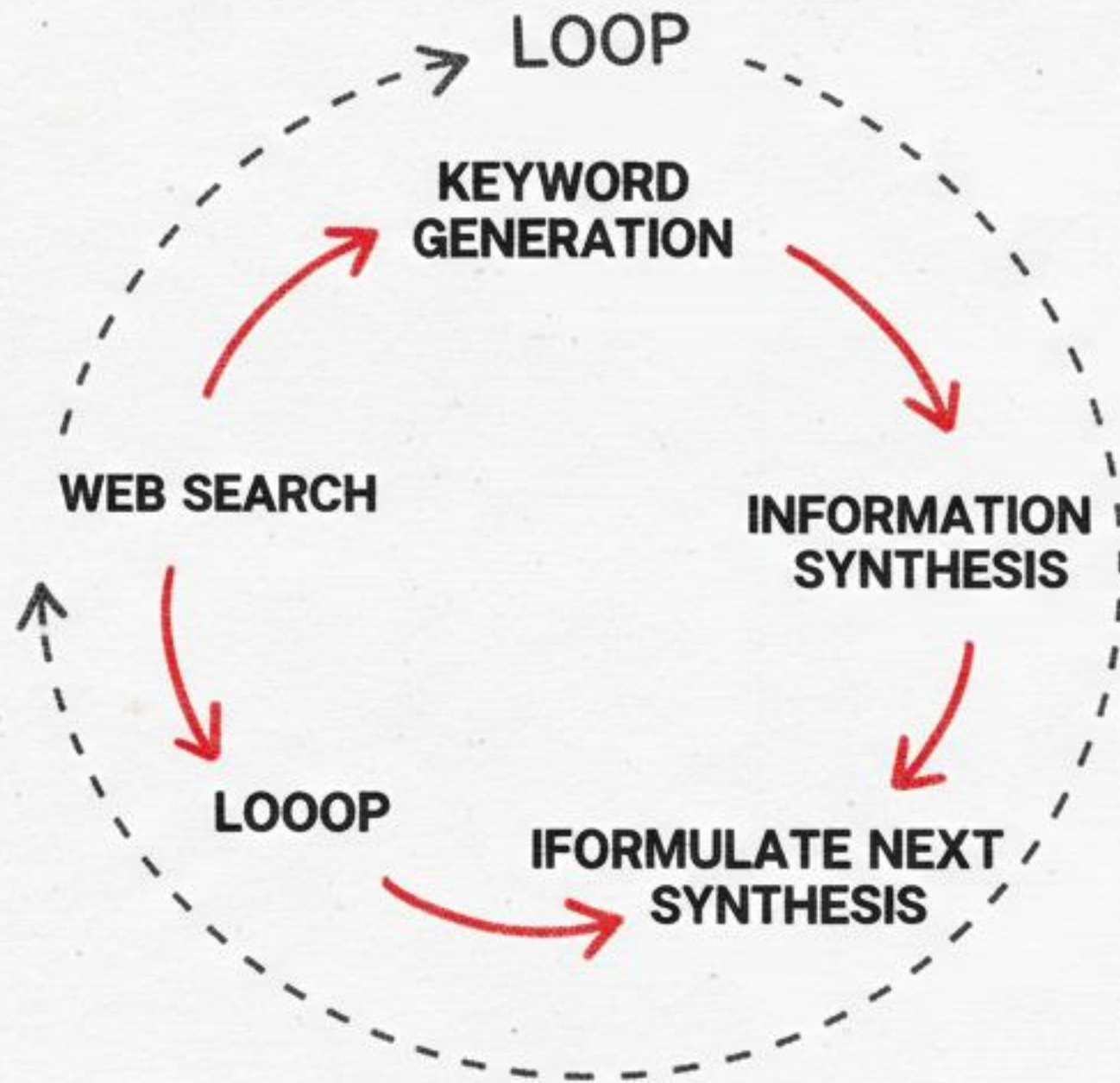Agents Break Down Task

**4. REPORT GENERATION**

Input: O, G)
Process: Reporv/ FDV –Viiarts (G)
Output: Repow/ FDV → Code – Final
Final Multinododal Report

**GOAL:** Generate Multimodorshal Report LIKE R. KEY
**UNIVERSAL:** Agents Cross-model, Cross-topic

# STAGE 1: RESEARCHING

LOOP

KEYWORD GENERATION

WEB SEARCH

INFORMATION SYNTHESIS

LOOOP

IFORMULATE NEXT SYNTHESIS

- **OBJECTIVE:** *Given topic t, obtain comprehesive, upded, and cited 'learnings' L through multi-round retrieval & reasning.*

- *LLM generates keywords K based t for Web search search.*

- *Utilize K to search web pages P; analyze & extract information.*

- *Synthtize results into structured 'learnings L; generate next generate next research question q.*

- *Iterare through n_R rounds & refine understanding of t.*

- **FINAL OUPUT:** *Learnings L with key info & exteraJtations, providing knowledge base for planning & reporting.*

# OUTPUT: LEARNED FACTS (L)

# VISUALIZATION TEXTUALIZATION & FDV: A FRAMEWARK

## (A) Original Visualization

### UK City Traffic Volume

### US City Traffic Volume

Visualzation Visualization mage input.

## (B) Formal Description of Visualization

Layout Two strp fi plume xteri'bfs kin ditt.
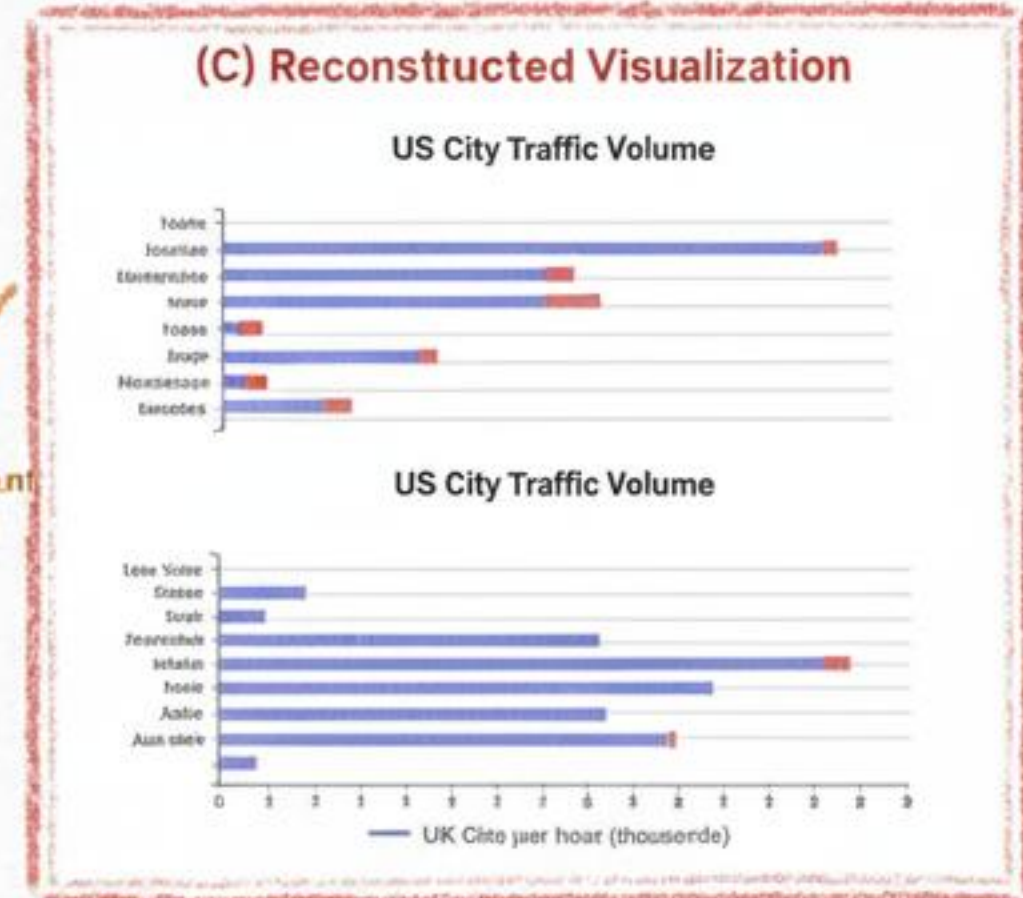- X tione 'frist vater arikemc ar tfvowling potwhp grisine
- @werernaied tes's prut tib1ee ia tie finheay wo ste ( US CM.

Lepout: Two strip tick mark 'nmorvis tmske sofics sa tie d tisne
- Tht ouist tilde Toing thvs st er tial aith of belove atsunt tien tik cinty atof.
- Ynoke : Categoric atcaltio vitnoe er he varotrrry 'fU S veliepvalle boxtle ix a ua s It psrnrel ovixe).

Cete: Trob tick motse noark tronce sa finr fuheg arc ttor FDK 'M woioh tindik Trck marikis cones ary; tte on a ertion ow tcc alt twoow tonie f FDf.

Marks : Srnoll fick owark't Irmeutuln horraitdimes trerted DG ?0 fe sioitus with tick mnume FDf,) onotor wohe ept is uris viora ffely.

Structured text description (FDV).

## (C) Reconstructed Visualization

### US City Traffic Volume

### US City Traffic Volume

Visualization generated from FDV

Extract Design

Implment Design

## Textualization Algorithm Key Steps:

1. FOR each chart ickarrt i in Report R:
2. FDV$i$ = M$v$) (Rnv(Image)   // Multimodal LLM extracts >R
3. Replace Image = Code_from FDV_R) // Reconstutrct for validation
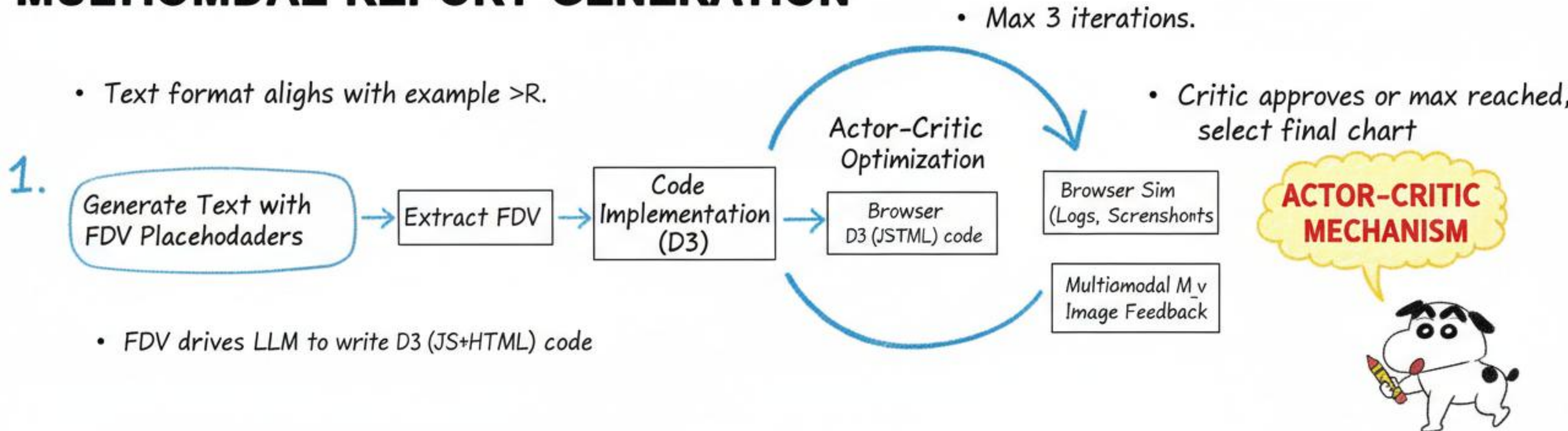4. END FOR_Pure text report >R

# STAGE 3 & 4: PLANNING & MULTIMODAL REPORT GENERATION

## PLANNING

- Based learnings L from multi-turn retrieval, topic t, and genand text examples >R. Visual Style Guide G
- Outline O: Hierarairical chapter structure with titles & summaries. Determines Determines narrative flow.
- Style Guide G: Learns color schemes, font hieraeicies, chart layouts from example reports. Ensures visual consistency.

## MULTIOMDAL REPORT GENERATION

- Max 3 iterations.

- Text format alighs with example >R.

- Critic approves or max reached, select final chart

Actor-Critic Optimization

1.

Generate Text with FDV Placehodaders → Extract FDV → Code Implementation (D3) → Browser D3 (JSTML) code

Browser Sim (Logs, Screenshonts

Multiomodal M_v Image Feedback

**ACTOR-CRITIC MECHANISM**

- FDV drives LLM to write D3 (JS+HTML) code

# EXPERIMENT & EVALUATION:
## MultimodalReportBench & Results Overview

- **DATA CONSTRUCTION:** Introducing MultimodalReportBench for systematic of multiodal report generation quality.

- Dataset contains 100 real-world topics, sourced from public multimodal report websites, authored by human experts.

- **DESIGNED** 5 dedicated evaluation metrics: Content Complleinss Quality, Text-Chart Style Consistenty, etc.

- **BASELINES:** Modifized DataNarative framework to generate chart pachsceraberk Visuslizatienss DataNarrative framework to generate chart placloores & text for ospen-source task.

- **EVALUATION METHOD:** Combining automatic & human evaluation to compare various compare propgeraity & open-source models.

- **KEY RESULT:** Using the same Claude 3.7 Sonnet model as generator, Multimodal achieved 82. 3% WIN RATE against baseline in overall assessment.

- Results show: FDV & Agent-based phased framework SIGNIFICANTY IMPROVE & & utility of interbabed text-chart reports.

82%

18%

Overall Win Rate: Multiomodal DeepResacher vs. Baseline

**METHOD SHOWS SUPEPHORITY IN MULTIPLE METRICS & HUMAN EVALUATION!**

# SUMMARY & OUTLOOK

## WORK SUMMARY

- Proposed novel task & benchwark for zero-shot zext-chot-chart intwhinined mutimodal report generation.

- Introduced FDV: a general, structured text represtation for in-conxt leand learning & automatic chart recorustrction.

- Designed Multimodal DeepResearcher 4-stage agent framework: integrating few-shot learning, planning planning & generation.

- Experiments show framework significanly outforoms Dataative baseline, achieving 82% overall win-rate with same model.

## FUTURE DIRECTIONS

- Expand to more visusization types (e.g, interactivs, animions) & complex forms, video, audio.

- Explore end-eed RL or self-play mecanisms to enhance enhance synragy of retireval, planning foms, & generation.

- Research finer multiodoal evaluation metrics & human-computer co-editing worfflows for practical deployment

## TAKE-HOME MESSAGE

**Structured Visual Descriptions + Stagent Colalboration = Significanty Improved LLM Performance in Real-World Mutlmodal Report Generation!**

# ACKOWLEGEMENTS

THANKS FOR THE SUPPORT AND DISCUSSIONS FROM MY
ADVISORS, LABMATES, AND COLLABARATORS.

APPRECIATION FOR THE CONTRIBUTORS OF PUBLIC DATASETS AND
OPEN-SOURCE TOOLS.

GRATEFUL FOR THE VALUABLE TIME AND FEEDBACK FROM REWIENCE MEMBERS.

# THANK YOU ALL!