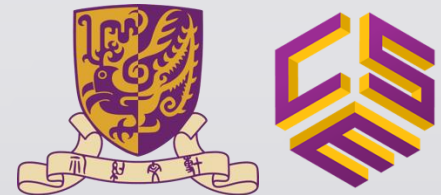# Scenethesis:
# Structure-Based XR Scene Synthesis

LYU2406

Supervisor: Professor Michael R. Lyu
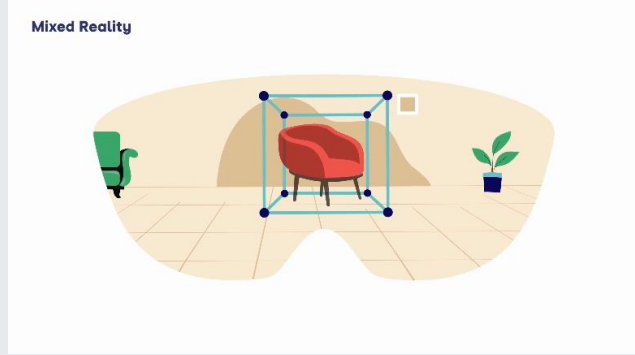
Presenter: LAM Yiu Fung Anson

# Outline

1. Background
2. Motivation
3. Methodology
4. Demo
5. Experiments
6. Conclusion
7. Future Work

# Background

| Extended Reality (XR) | | |
|---|---|---|
| **Virtual Reality (VR)** | **Augmented Reality (AR)** | **Mixed Reality (MR)** |
| Brings users to a <u>fully</u> simulated and isolated world | Combines the real and virtual worlds by <u>overlaying</u> digital content onto a dedicated device | Similar to AR, but things happening in the physical world can <u>affect</u> the virtual world |

# Background



Mobile augmented reality (AR) users worldwide from 2023 to 2028 (in millions)

# Motivation

| Software testing | |
|---|---|
| **Non-XR software** | **XR (AR and MR) software** |
| 1. Use DFS to go through every function<br>2. Perform unit tests automatically<br><br>😍😍😍 | 1. Design and construct various physical environments<br>2. Test all functions within each environment<br>3. Analyze the recordings manually to check for guideline violation<br><br>🥵🥵🥵 |

Term 1

**Can we also automate XR testing?**

# Motivation (by comparing related work)

| Non-playable scene | Playable scene |
|---|---|
| A unified mesh → not physically correct | Interactable with individual objects |



Text2Room, ICCV 2023



InstructScene, ICCV 2024



BlockFusion, SIGGRAPH 2024



Holodeck, CVPR 2024

# **Motivation** (by comparing related work)

- Scene graph: node as object, edge as relation



➤ Categorical nodes, $o \in \{1, ..., K_o\}$

➤ Categorical edges, $r \in \{1, ..., K_r\}$

$K_o, K_r$ are predefined!

1. Scene-specific, not <u>flexible</u> and practical for XR testing

2. Cannot <u>accurately</u> describe all relations

Image source: InstructScene

# **Motivation** (by comparing related work)

We need a structure that is:

1. Human-understandable: for explainability and modifiability

2. Unambiguous: all (physically) possible relations in a scene can be modeled

3. Flexible: able to generate a diverse set of scenes

# Methodology – ScenethesisLang

- Domain-specific language (DSL)

- Describes diverse, realistic, and physically plausible 3D scenes

- Focuses on expressiveness, human-readability, generative capability, and physical plausibility

# Methodology – ScenethesisLang

$scene \in Scenes ::=$ **SceneType:** $scene\_type$;

**SceneDescription:** $scene\_description$;

**Regions:** $regions$;

**Connections:** $connections$;

**Constraints:** $constraints$;

**ProbabilisticParameters:** $probabilistic\_parameters$;

**TemporalRequirements:** $temporal\_requirements$

LYU2406

# Methodology – ScenethesisLang

- Scene type: overall nature of the environment (indoor/outdoor)

- Scene description: textual description of the scene

**SceneType:** *scene_type*;

**SceneDescription:** *scene_description*;

**Regions:** *regions*;

**Connections:** *connections*;

**Constraints:** *constraints*;

**ProbabilisticParameters:** *probabilistic_parameters*;

**TemporalRequirements:** *temporal_requirements*

# Methodology – ScenethesisLang

- Region: spatial subdivision within a scene (e.g., room, outdoor area)
  - An outdoor area is defined as the bottom shape (from bird's-eye view) of a prism with infinite height.

**SceneType:** *scene_type*;

**SceneDescription:** *scene_description*;

**Regions:** *regions*;

**Connections:** *connections*;

**Constraints:** *constraints*;

**ProbabilisticParameters:** *probabilistic_parameters*;

**TemporalRequirements:** *temporal_requirements*

$regions \in \textit{Regions} ::= region \mid region;\ regions$

$region \in \textit{Region} ::= id \leftarrow \textbf{region(}description, shape, materials,$
$construction, objects, visibility, constraints\textbf{)}$

$shape \in \textit{Shape} ::= cuboid \mid ellipsoid \mid mesh \mid$
$shape \cup shape \mid shape \cap shape \mid shape - shape$

$cuboid \in \textit{Cuboid} ::= \textbf{cuboid(}corner.min, corner.max\textbf{)}$

$ellipsoid \in \textit{Ellipsoid} ::= \textbf{ellipsoid(}center, a, b, c\textbf{)}$

$mesh \in \textit{Mesh} ::= \textbf{mesh(}vertices, faces\textbf{)}$

$corner \in \textit{Corner} ::= \textbf{min} \mid \textbf{max}$

$materials \in \textit{Materials} ::= floor\_material;\ non\_floor\_material$

$floor\_material \in \textit{Material} ::= \textbf{material(}color, texture, reflectivity, opacity\textbf{)}$

$non\_floor\_material \in \textit{Material} ::= \textbf{material(}color, texture, reflectivity, opacity\textbf{)}$

$construction \in \textit{Construction} ::= vertices, faces$

$vertices \in \textit{Vertices} ::= \textbf{[}(x, y, z), ...\textbf{]} \mid \textbf{[}(x, y), ...\textbf{]}$

$faces \in \textit{Faces} ::= \textbf{[}(index1, index2, index3), ...\textbf{]}$

Wall? Ceiling?

# Methodology – ScenethesisLang

- Object: entity within a region

$$objects \in Objects ::= object \mid object;\ objects$$

$$object \in Object ::= id \leftarrow \textbf{object(}category, description, dimensions,$$
$$position, rotation, bounding\_box,$$
$$probabilistic\_properties, visibility\textbf{)}$$

$$category \in Categories ::= \textbf{text}$$

$$description \in ObjectDescription ::= text$$

$$dimensions \in Dimensions ::= \textbf{(length, width, height)}$$

$$position \in Position ::= \textbf{(x, y, z)} \mid \textbf{offset(position, vector)} \mid$$
$$\textbf{relativeTo(reference)}$$

$$rotation \in Rotation ::= \textbf{(roll, pitch, yaw)} \mid \textbf{relativeTo(reference)}$$

$$bounding\_box \in BoundingBox ::= \textbf{bounding\_box(}min, max\textbf{)}$$

$$visibility \in Visibility ::= \textbf{rayTrace(}density, occlusion\textbf{)} \mid$$
$$\textbf{regionVisibility(}region, conditions\textbf{)}$$

---

**SceneType:** *scene_type*;

**SceneDescription:** *scene_description*;

**Regions:** *regions*;

**Connections:** *connections*;

**Constraints:** *constraints*;

**ProbabilisticParameters:** *probabilistic_parameters*;

**TemporalRequirements:** *temporal_requirements*

# Methodology – ScenethesisLang

- Object: entity within a region

$$lights \in Lights ::= light \mid light;\ lights$$
$$light \in Light ::= id \leftarrow \textbf{light}(category, description, intensity, position, visibility)$$
$$intensity \in Intensity ::= \textbf{float} \mid \textbf{distribution}$$

**SceneType:** $scene\_type$;

**SceneDescription:** $scene\_description$;

**Regions:** $regions$;

**Connections:** $connections$;

**Constraints:** $constraints$;

**ProbabilisticParameters:** $probabilistic\_parameters$;

**TemporalRequirements:** $temporal\_requirements$

# Methodology – ScenethesisLang

- Connection: spatial relationship between two regions

$$connections \in Connections ::= connection \mid connection; \ connections$$
$$connection \in Connection ::= \textbf{connection}(region1, region2, type,$$
$$description\textbf{)}$$

**SceneType:** $scene\_type;$

**SceneDescription:** $scene\_description;$

**Regions:** $regions;$

**Connections:** $connections;$

**Constraints:** $constraints;$

**ProbabilisticParameters:** $probabilistic\_parameters;$

**TemporalRequirements:** $temporal\_requirements$

# Methodology – ScenethesisLang

- Constraint: requirement that ensures physical plausibility or reasonableness, or that meets user-defined specifications

SceneType: *scene_type*;

SceneDescription: *scene_description*;

Regions: *regions*;

Connections: *connections*;

Constraints: *constraints*;

ProbabilisticParameters: *probabilistic_parameters*;

TemporalRequirements: *temporal_requirements*

$constraints \in Constraints ::= constraint \mid constraint; constraints$
$constraint \in Constraint ::= spatial\_condition \mid probabilistic\_condition \mid$
$temporal\_condition \mid object\_relation \mid$
$visibility\_condition \mid physical\_constraint \mid$
$user\_defined(logic, priority)$

$spatial\_condition \in SpatialCondition ::=$ **inside(region)** | **outside(region)** |
**above(object, height)** |
**below(object, height)** |
**nearby(object, distance)** |
**alignedWith(object, axis)** |
**tangentTo(surface)** |
**distanceBetween(object1, object2) == d**

$visibility\_condition \in VisibilityCondition ::=$ **canSee(observer, target)** |
**occludes(object1, object2)** |
**visibleInRegion(observer, region)** |
**rayTraceValid(observer, target, density)**
$physical\_constraint \in PhysicalConstraint ::=$ **noCollision(object1, object2)** |
**collisionFreeRegion(region)** |
**stablePosition(object)** |
**gravityAligned(object)**
$user\_defined \in UserDefinedConstraint ::=$ **customLogic(logicExpression, priority)**

$object\_relation \in ObjectRelation ::=$ **relation(object1, object2, relation_type)**
$relation\_type \in RelationTypes ::=$ **above** | **below** | **inside** | **outside** | **nearby** | **aligned** |
**occludes** | **intersects**

# Methodology – ScenethesisLang

- Probabilistic parameter: ranges or distributions for attributes such as object positions

$probabilistic\_condition \in ProbabilisticCondition ::=$

**probability(p): condition | distributionBased(object, param: distribution)**

**SceneType:** $scene\_type$;

**SceneDescription:** $scene\_description$;
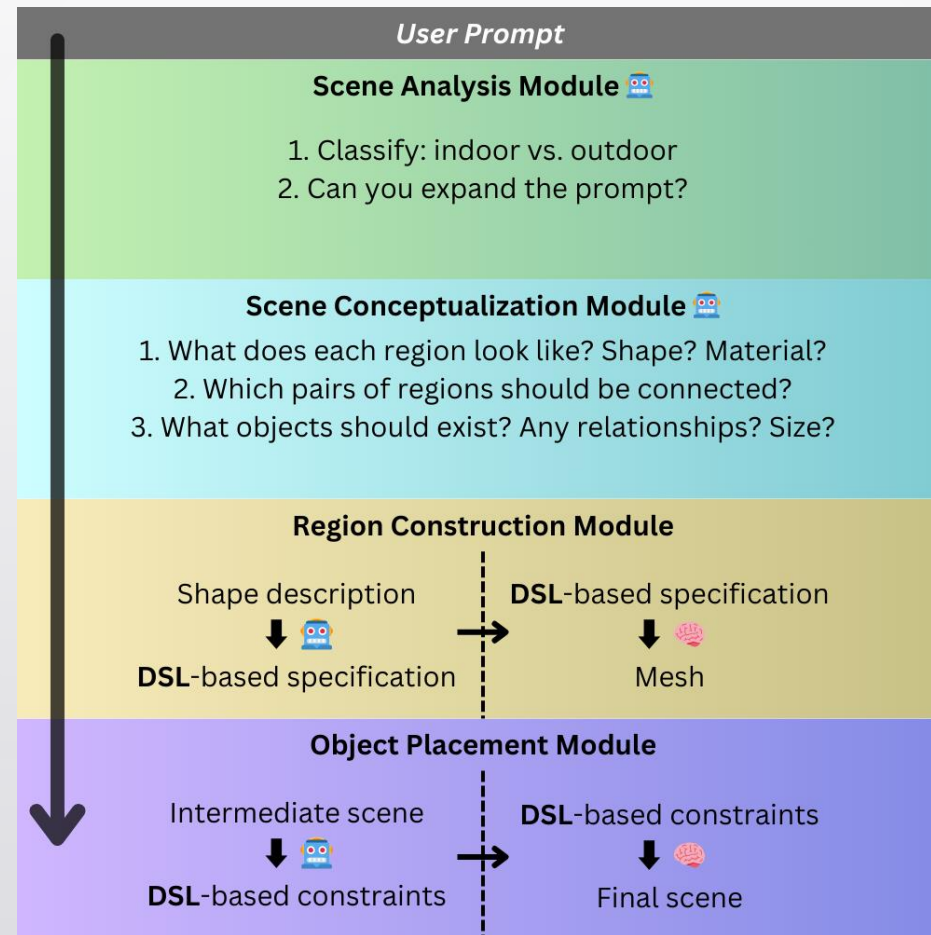
**Regions:** $regions$;

**Connections:** $connections$;

**Constraints:** $constraints$;

**ProbabilisticParameters:** $probabilistic\_parameters$;
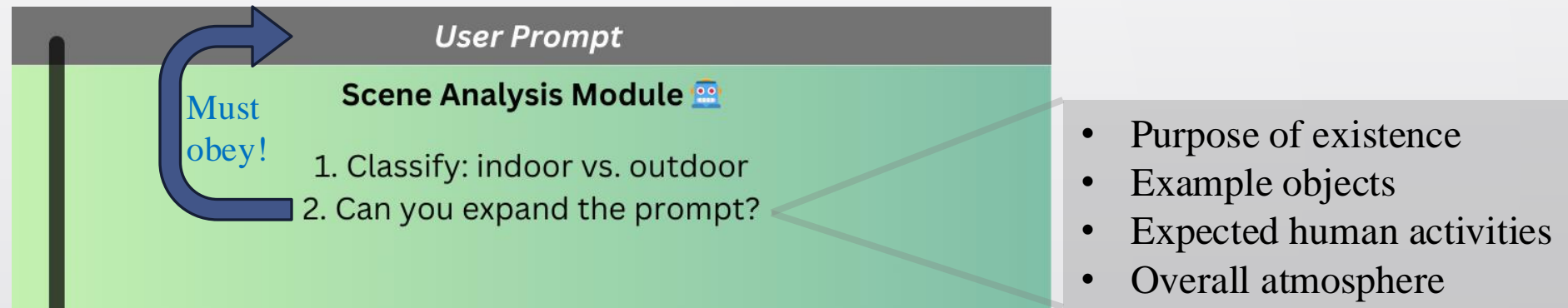
**TemporalRequirements:** $temporal\_requirements$

# Methodology – Scenethesis



🤖 : LLM

# Methodology – Scenethesis

- Scene Analysis Module: enhances the overall understanding of the desired scene



**User Prompt**

**Scene Analysis Module** 🤖

Must obey!

1. Classify: indoor vs. outdoor
2. Can you expand the prompt?

- Purpose of existence
- Example objects
- Expected human activities
- Overall atmosphere

$scene \in Scenes ::=$ **SceneType:** $scene\_type$;

**SceneDescription:** $scene\_description$;

# Methodology – Scenethesis

- Scene Conceptualization Module: creates a semantic draft of the desired scene

  - Semantic: only <u>text</u> is outputted

- Purpose of existence
- Example objects
- Expected human activities
- Overall atmosphere

3 prompts per region:
1. Names of all objects
2. Relationships
3. Attributes

**Scene Conceptualization Module** 🤖
1. What does each region look like? Shape? Material?
2. Which pairs of regions should be connected?
3. What objects should exist? Any relationships? Size?

$$m^* = \underset{m \in \mathcal{D}_m}{\arg\max} \, \mathrm{CLIP}\,(m, q_m)$$

$$(o^*, d^*) = \underset{(o,d) \in \mathcal{D}_o}{\arg\max} \, \frac{\alpha \cdot \underset{r \in \mathcal{R}}{\max} \, \mathrm{CLIP}\,(r, q_o) + \beta \cdot \mathrm{SBERT}\,(d, q_o)}{\alpha + \beta}$$

$q_o$: A 3D object of a {category} named {name} ({description})

$region \in Region ::= id \leftarrow \textbf{region(}description, shape, materials,$
$construction, objects, visibility, constraints\textbf{)}$

$connection \in Connection ::= \textbf{connection(}region1, region2, type,$
$description\textbf{)}$

$object \in Object ::= id \leftarrow \textbf{object(}category, description, dimensions,$
$position, rotation, bounding\_box,$
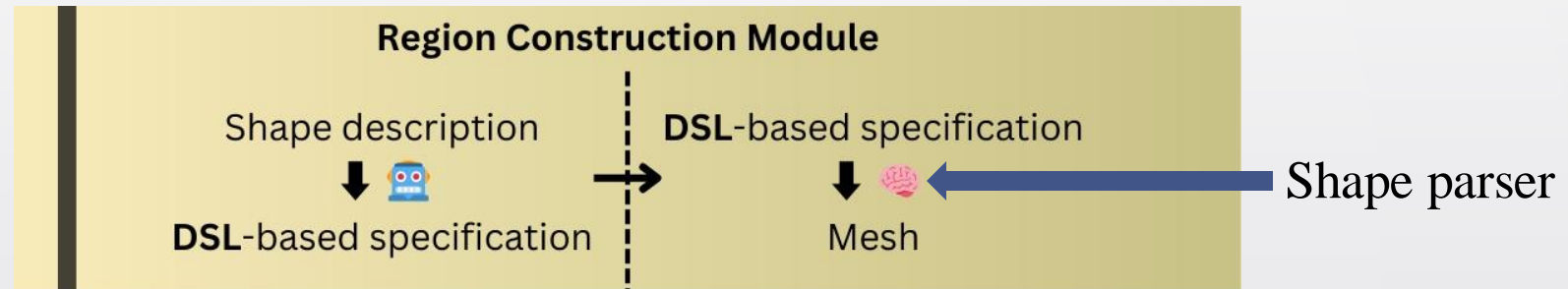$probabilistic\_properties, visibility\textbf{)}$

# Methodology – Scenethesis

- Region Construction Module: constructs the actual boundaries of each region



Shape parser

$$shape \in \textit{Shape} ::= cuboid \mid ellipsoid \mid mesh \mid$$
$$shape \cup shape \mid shape \cap shape \mid shape - shape$$
$$cuboid \in \textit{Cuboid} ::= \mathbf{cuboid(}corner.min, corner.max\mathbf{)}$$
$$ellipsoid \in \textit{Ellipsoid} ::= \mathbf{ellipsoid(}center, a, b, c\mathbf{)}$$
$$mesh \in \textit{Mesh} ::= \mathbf{mesh(}vertices, faces\mathbf{)}$$
$$corner \in \textit{Corner} ::= \mathbf{min} \mid \mathbf{max}$$

$$construction \in \textit{Construction} ::= vertices, faces$$
$$vertices \in \textit{Vertices} ::= \mathbf{[}(x, y, z), ...\mathbf{]} \mid \mathbf{[}(x, y), ...\mathbf{]}$$
$$faces \in \textit{Faces} ::= \mathbf{[}(index1, index2, index3), ...\mathbf{]}$$

# Methodology – Scenethesis

- Object Placement Module: places selected objects to their optimal position with optimal orientation



Constraint solver

# Demo

```
genxr "A bedroom connected to a living room"
```

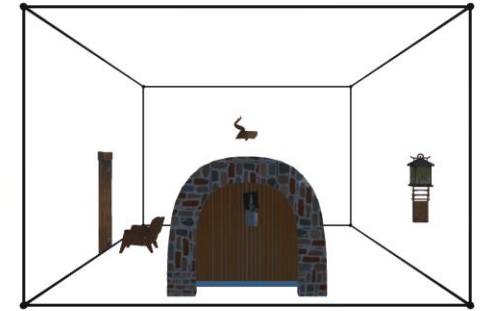User prompt:
"A bedroom connected to a living room"

# Experiments – Implementation

- Materials and objects are from Holodeck

- $\alpha = 100, \beta = 1$

- LLM: gpt-4o-2024-08-06

- Temperature: 0

- Machine: MacBook Pro (M1 Pro CPU, 16 GB RAM)
  - Mean execution time (before launching Blender and Unity): 2.5 minutes

# Experiments – Evaluation Metrics

- Quantitative:
  - $\text{Score}_{\text{CLIP}}(r, q) = (\text{CLIP}(r, q) + 1) \times 50$
    - $q_1$: "an image of a vibrant indoor scene"
    - $q_2$: user prompt
    - $q_3$: generated scene description
  - $\text{Score}_{\text{SBERT}}(d, p) = (\text{SBERT}(d, p) + 1) \times 50$
    - $d$ is the scene description, $p$ is the user prompt

Wireframe rendering $r$:
Pitch $\in \{0, 30, 45, 60, 90\}$
Yaw $\in \{0, 15, \dots, 330, 345\}$

# Experiments – Evaluation Metrics



- Qualitative:
  - We give the bird's-eye view to GPT-4o and ask:
    1. Does the generated scene contain every region mentioned in the prompt?
    2. Does the generated scene contain every object mentioned in the prompt?
    3. Is the generated scene physically plausible (e.g., are there any objects colliding with region boundaries or other objects)?
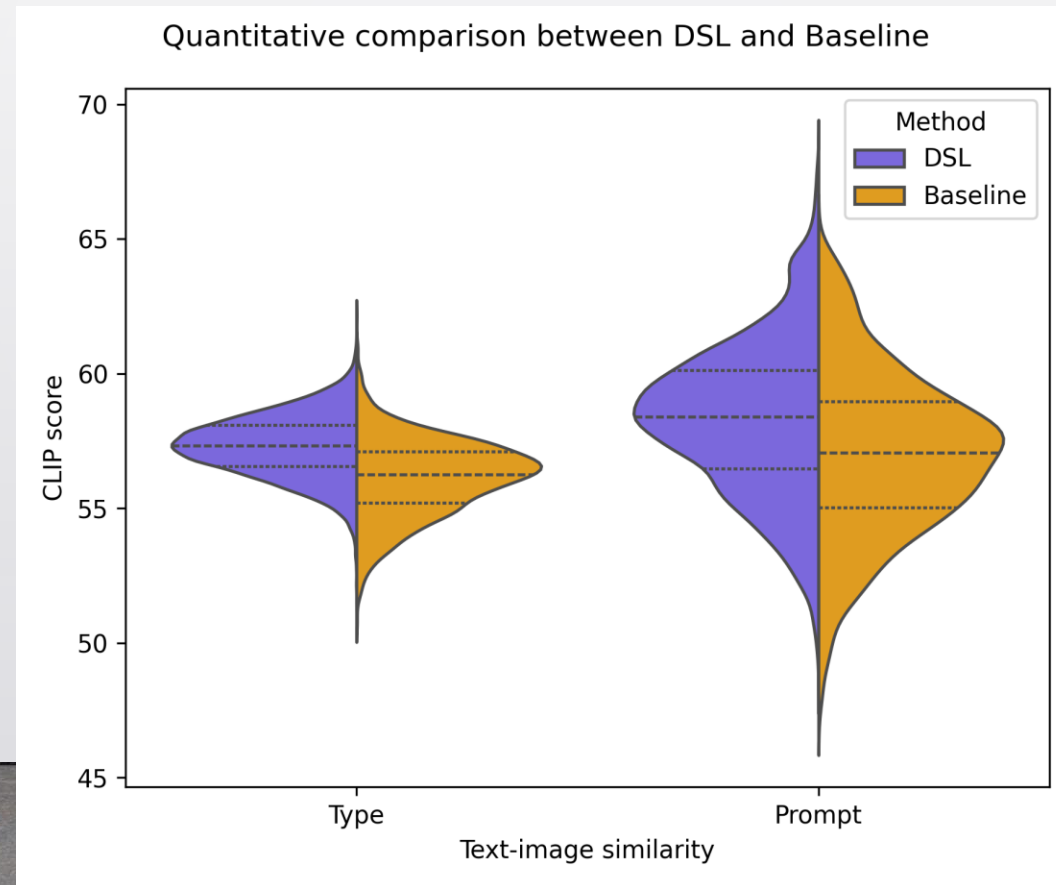    4. Is the generated scene visually pleasing?
  - Then rate from 1 to 10 with explanation

# Experiments – Results

1. Comparing with baseline

   - Baseline: Non-DSL injected to obtain a minimal output

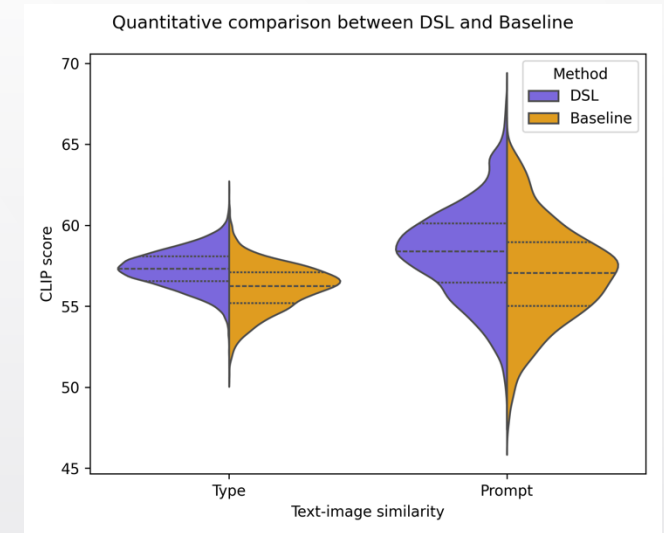   - 50 generated prompts, 3 trials per prompt

# Experiments – Results

1. Comparing with baseline

   - Quantitative results:



Quantitative comparison between DSL and Baseline

# Experiments – Results

1. Comparing with baseline

   - Quantitative results:

     - Small differences in CLIP scores

     - BUT! CLIP models were trained with real-life images, many of which have a main subject

     - Insensitiveness + small differences → Perhaps significant



Quantitative comparison between DSL and Baseline

# Experiments – Results

1. Comparing with baseline

   - Quantitative results:

     - Scenethesis generate significantly more objects → Vibrant and realistic

|  | Average # of regions | Average # of objects |
|---|---|---|
| **Scenethesis** | **1.36** | **20.973** |
| **Baseline** | 1.35 | 3.63 |

# Experiments – Results

1. Comparing with baseline

   - Qualitative results:



Qualitative comparison between DSL and Baseline

# Experiments – Results



Qualitative comparison between DSL and Baseline

1. Comparing with baseline

   - Qualitative results:

     - Similar scores for "region" and "physical"

       - Scenethesis performs better in general

     - Significant performance gaps for "object" and "visual"

       - An evidence that Scenethesis can better follow user's requirements while still producing visually pleasing scenes

# Experiments – Results

1. Comparing with baseline

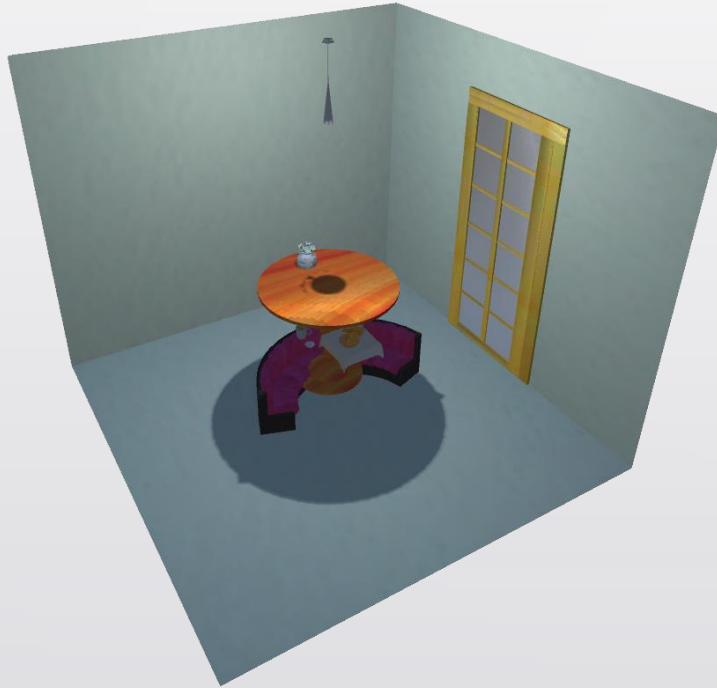   - Examples (left is Scenethesis, right is baseline):



*"A colorful nursery with a crib, rocking chair, and playful decor."*

# Experiments – Results

1. Comparing with baseline

   - Examples (left is Scenethesis, right is baseline):



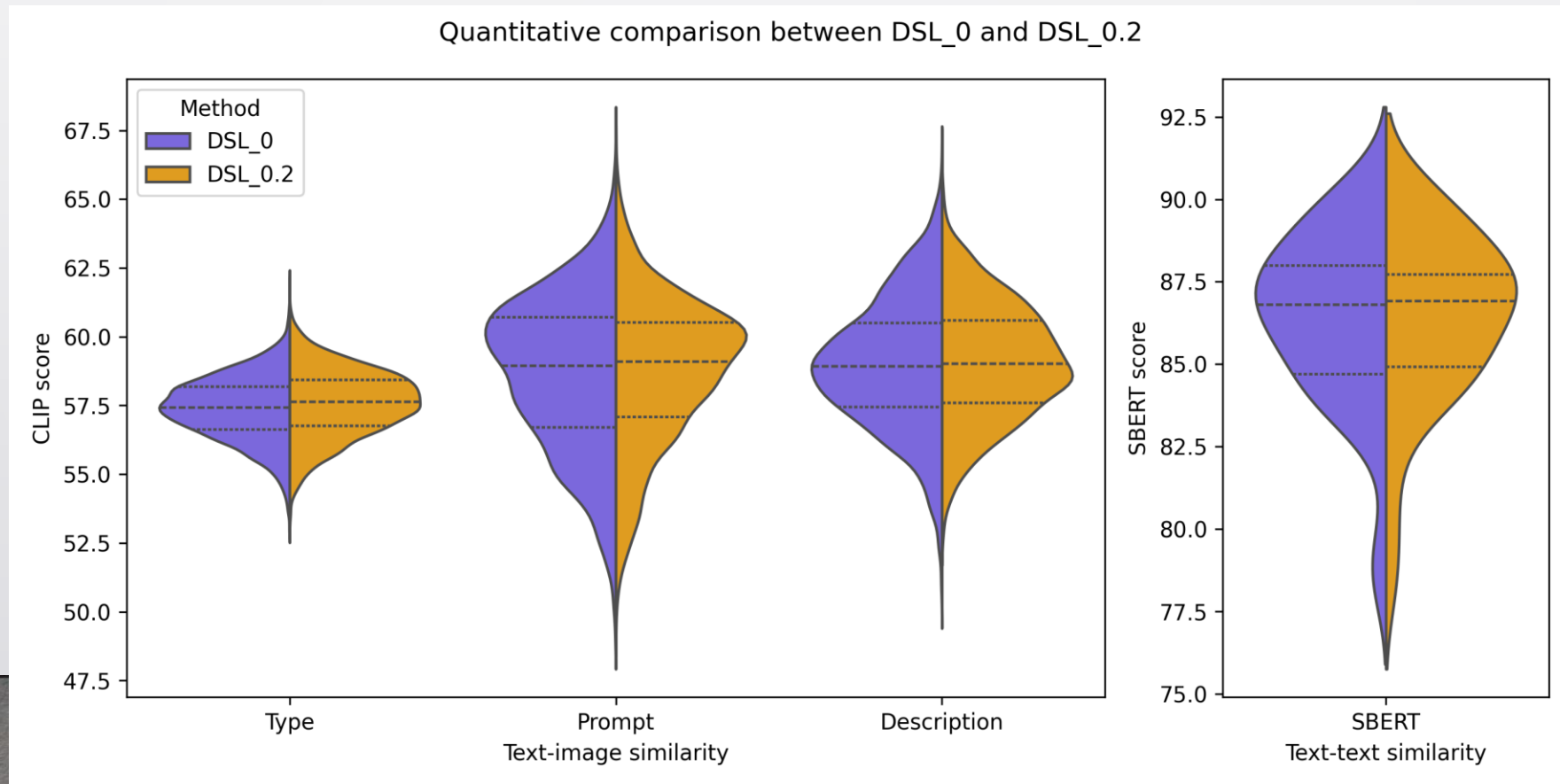*"A quaint breakfast nook with a round table and cushioned bench seating."*

# Experiments

2. Comparing different temperatures
   - $t \in \{0, 0.2, 0.5\}$
   - 20 generated prompts, 2 trials per prompt
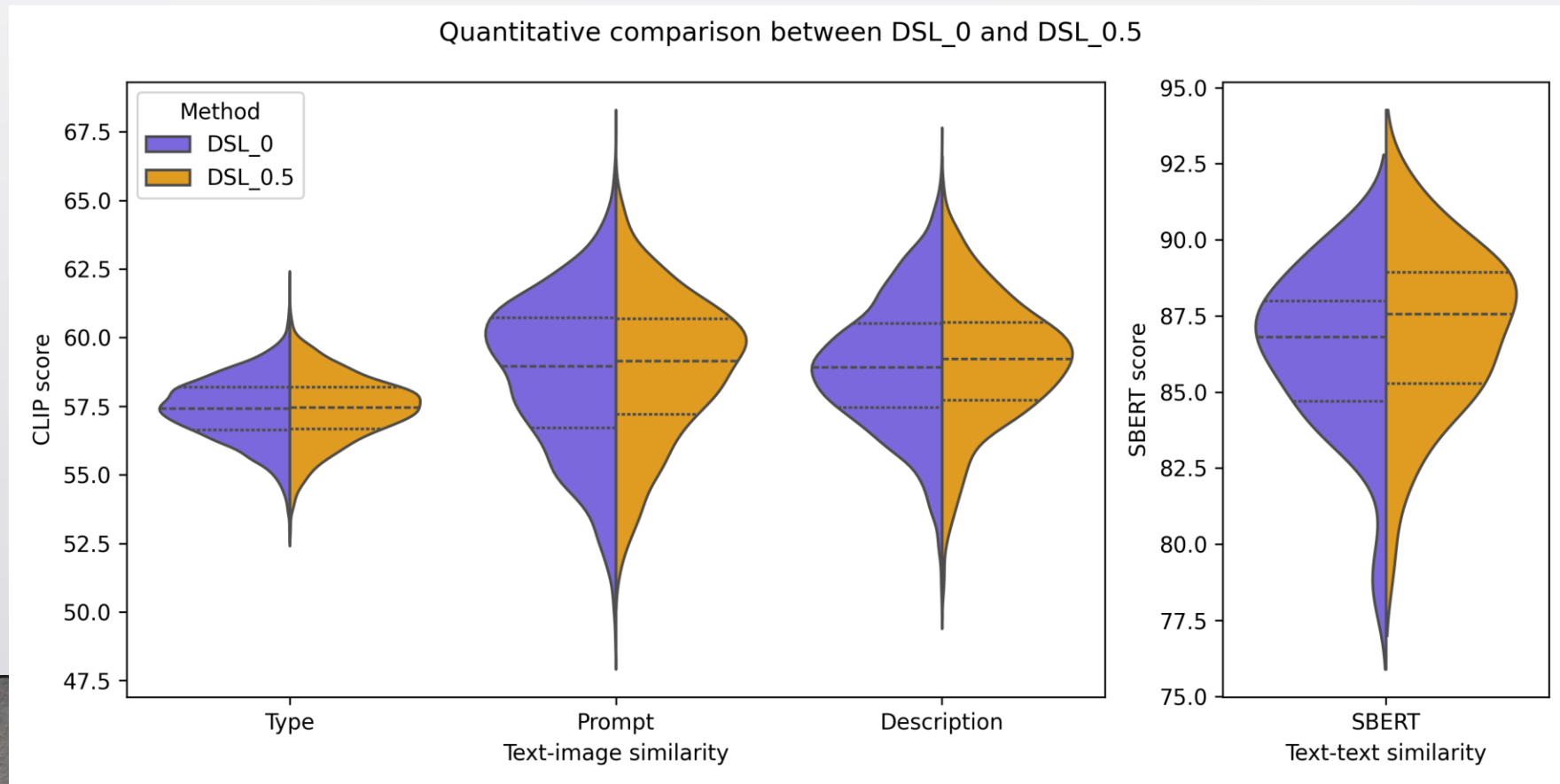
# Experiments

2. Comparing different temperatures

   - Quantitative results: 0 vs 0.2



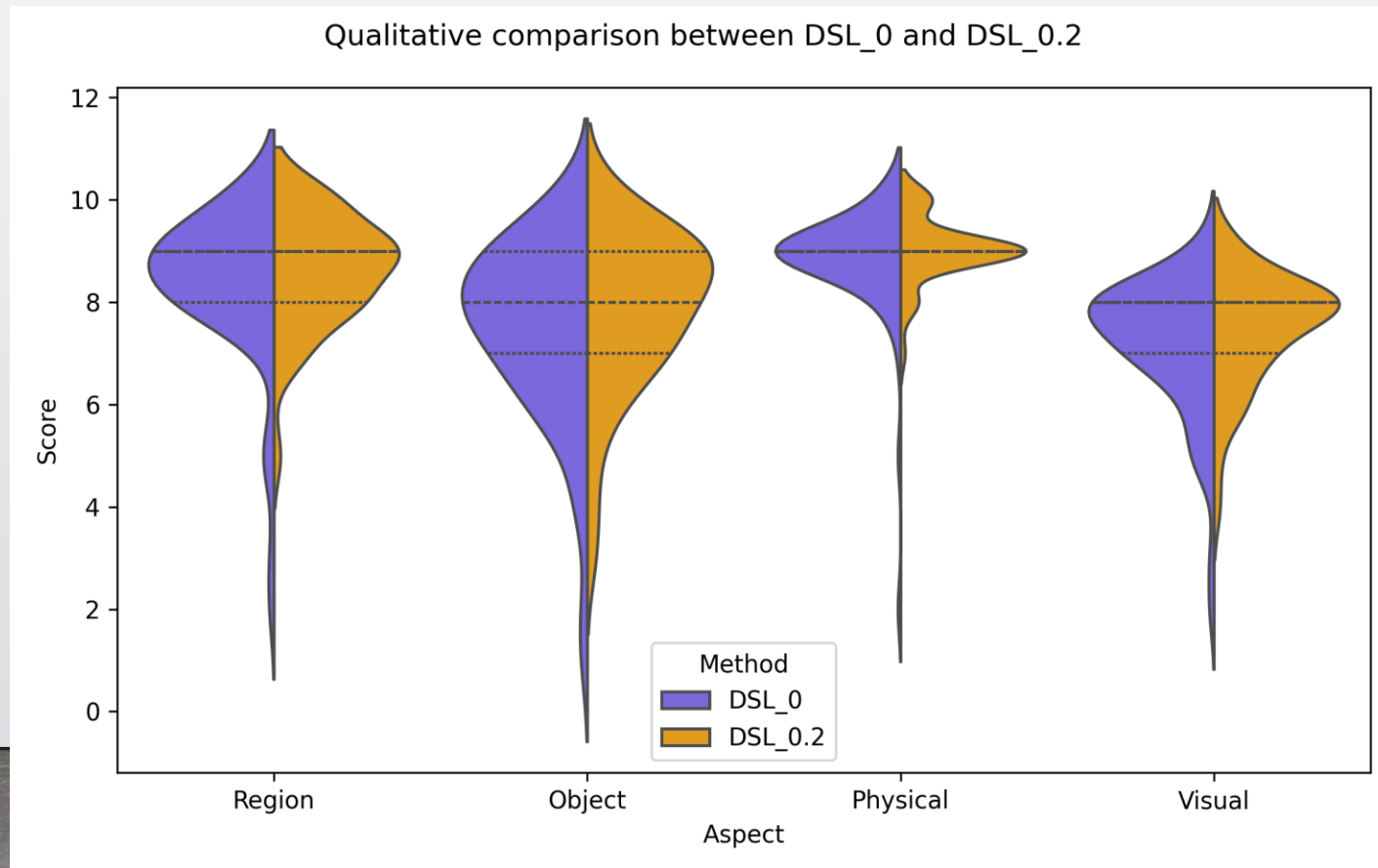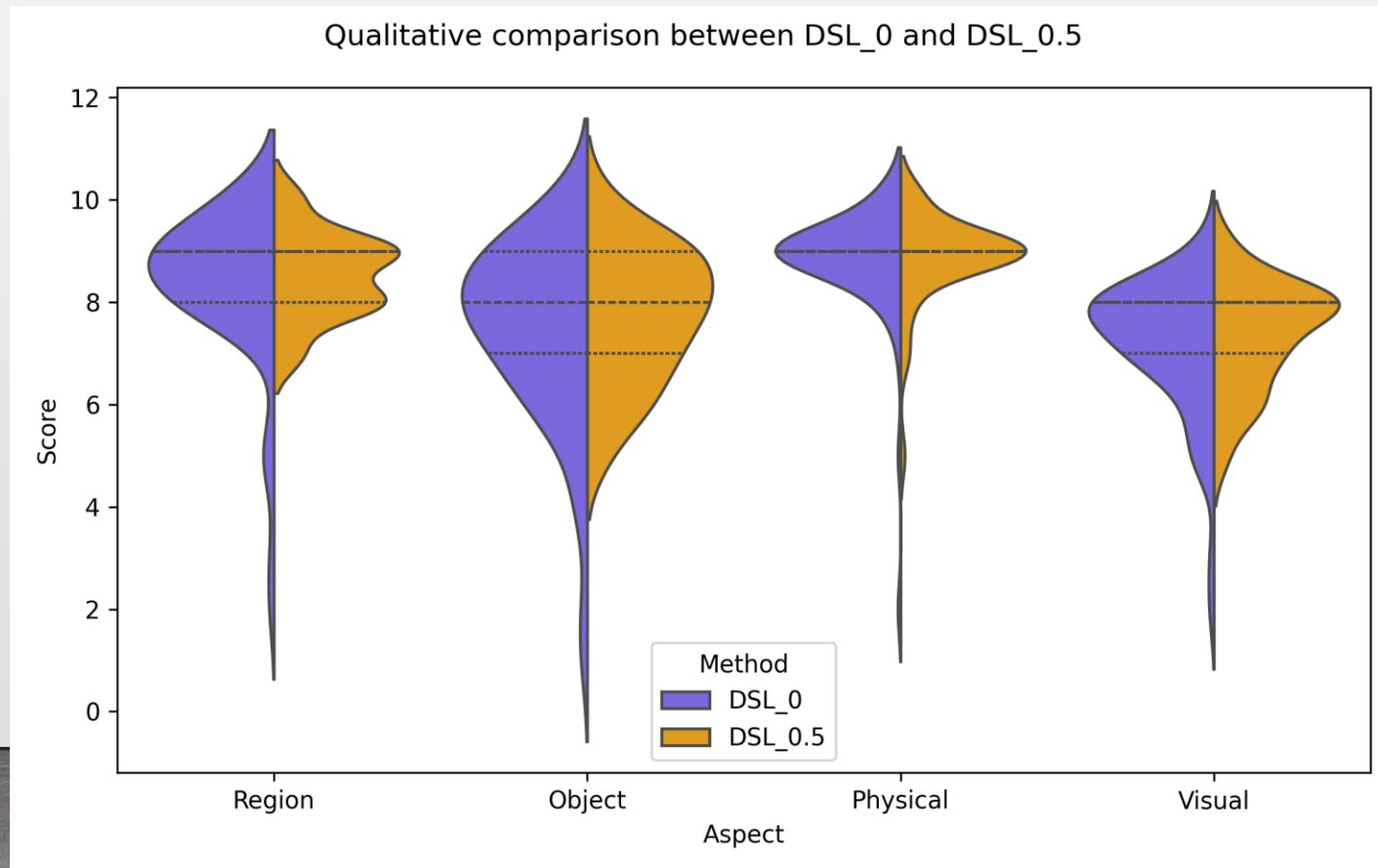Quantitative comparison between DSL_0 and DSL_0.2

# Experiments

2. Comparing different temperatures

   - Quantitative results: 0 vs 0.5

# Experiments

2. Comparing different temperatures

- Qualitative results: 0 vs 0.2



Qualitative comparison between DSL_0 and DSL_0.2

# Experiments

2. Comparing different temperatures

- Qualitative results: 0 vs 0.5



Qualitative comparison between DSL_0 and DSL_0.5

# Experiments

2. Comparing different temperatures

   - Subjective visual difference is not significant

   - So, temperature is not a major factor on the performance of Scenethesis
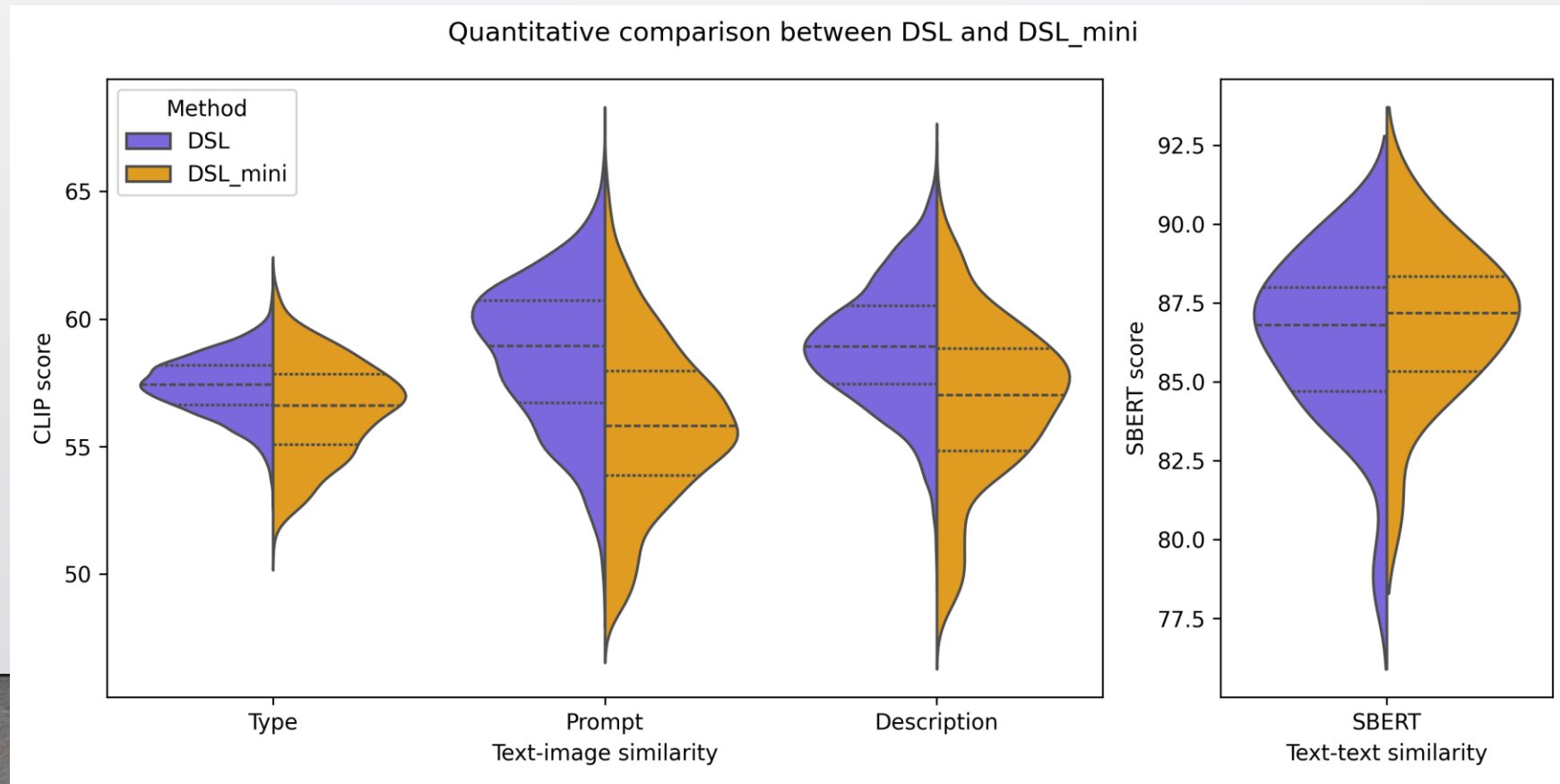
# Experiments

3. Comparing different LLMs

   - Model: gpt-4o-2024-08-06, gpt-4o-mini-2024-07-18, claude-3-5-sonnet-20241022
   - 20 generated prompts (same as before), 2 trials per prompt
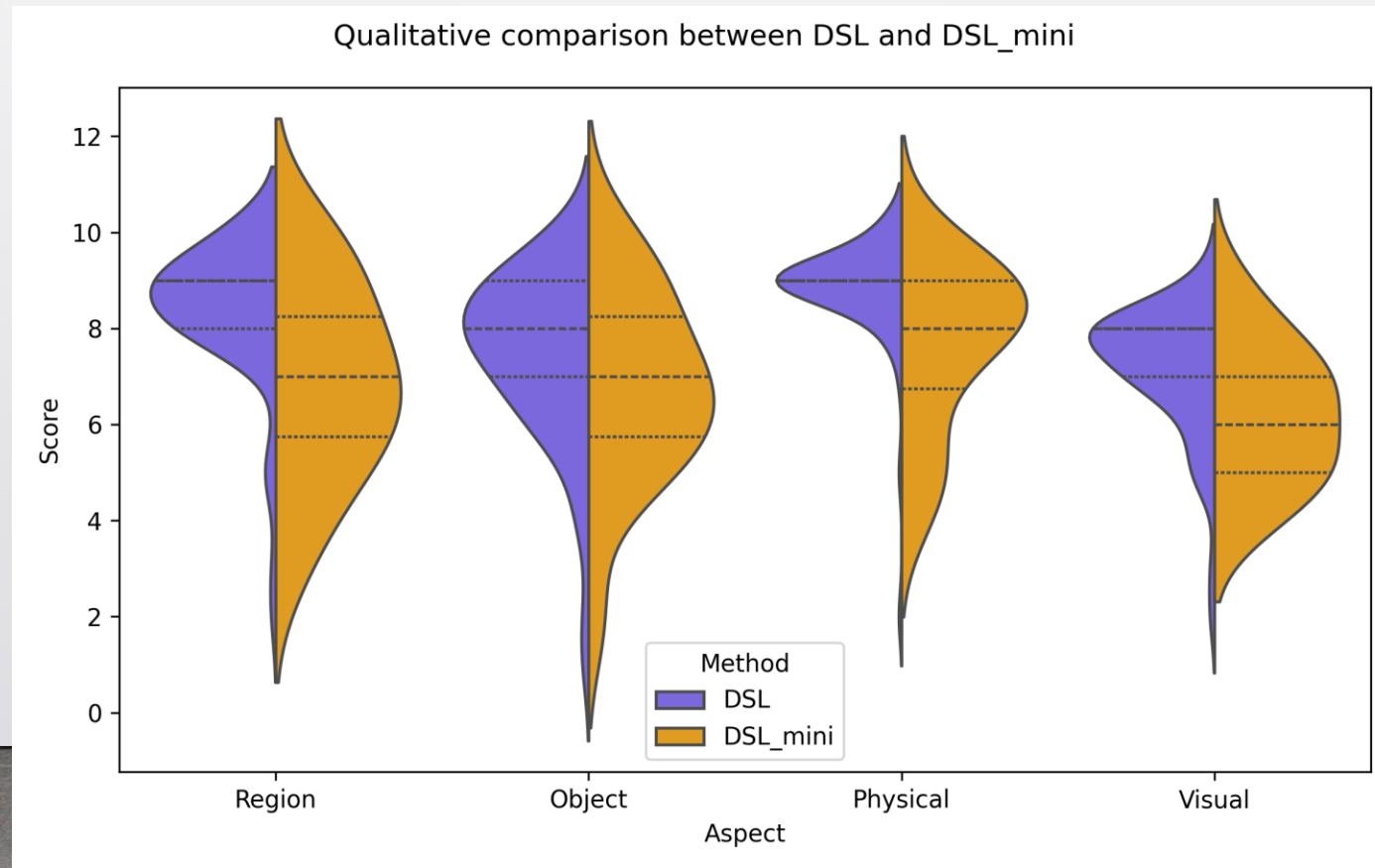
# Experiments

3. Comparing different LLMs

• Quantitative results: gpt-4o-mini-2024-07-18

# Experiments

3.  Comparing different LLMs

    - Qualitative results: gpt-4o-mini-2024-07-18



Qualitative comparison between DSL and DSL_mini

# Experiments
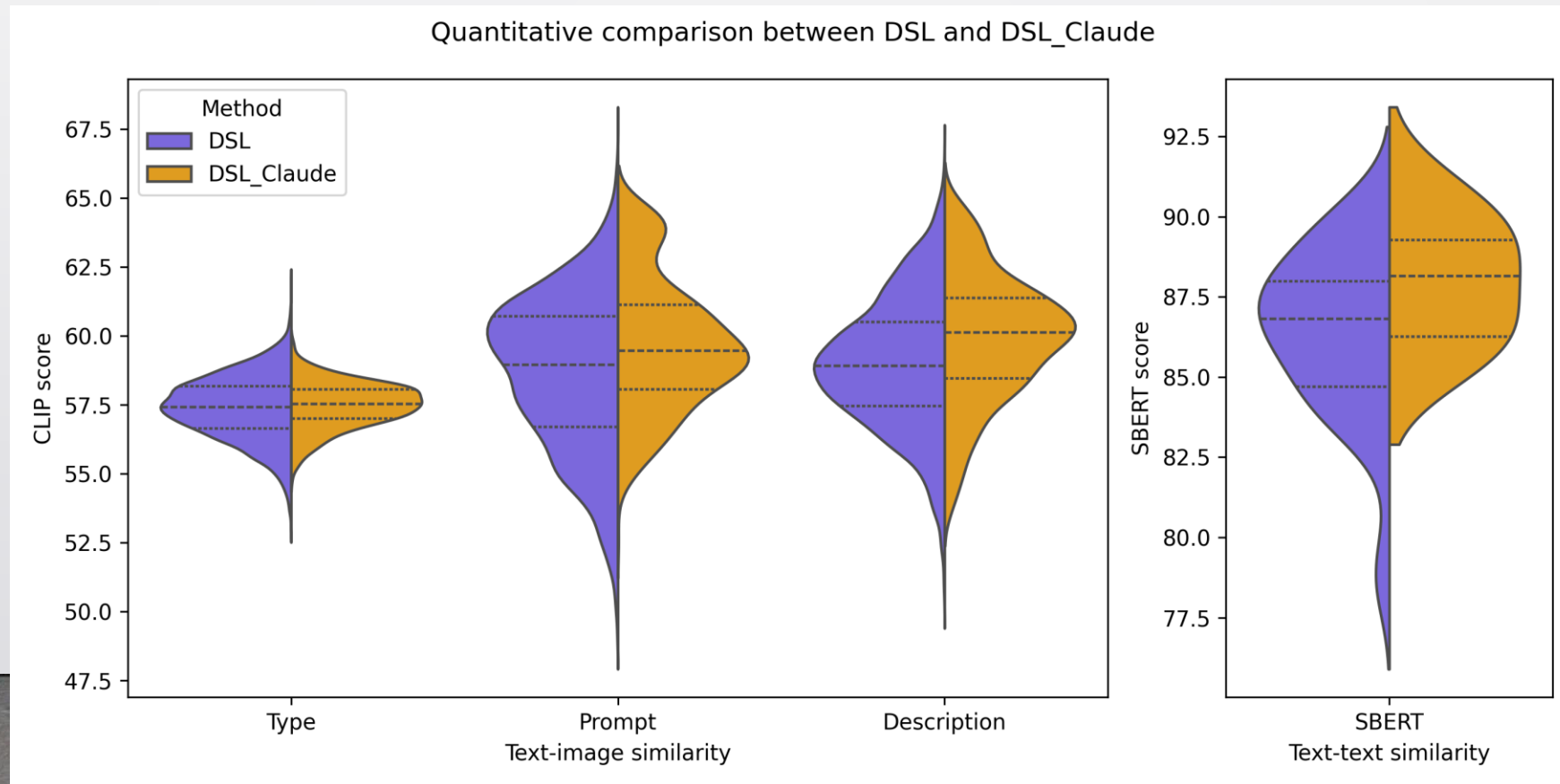
3. Comparing different LLMs

   - As expected, larger models with better processing and generative capabilities work better

# Experiments

3. Comparing different LLMs

- Quantitative results: claude-3-5-sonnet-20241022 (only 5 scenes could be generated)



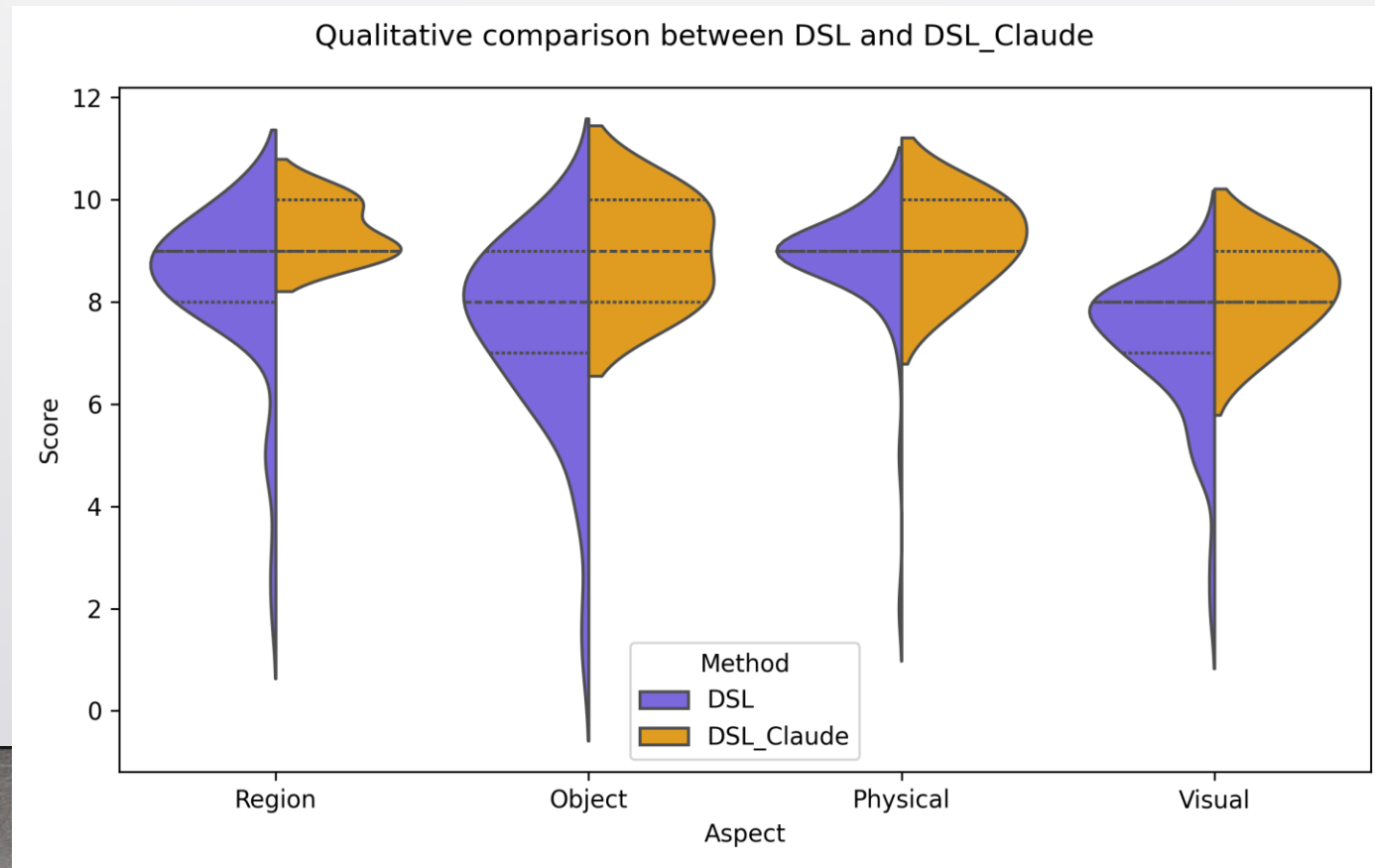Quantitative comparison between DSL and DSL_Claude

# Experiments

3. Comparing different LLMs

- Qualitative results: claude-3-5-sonnet-20241022 (only 5 scenes could be generated)

# Experiments

3. Comparing different LLMs
   - An inspiration: models that are better at writing → better scene description → better scene

# Conclusion

We propose:

1. ScenethesisLang, a novel domain-specific language designed for unambiguously describing a 3D scene; and

2. Scenethesis, a pipeline that takes a user prompt and generates a corresponding 3D scene using ScenethesisLang.

Experiments have demonstrated the potential of Scenethesis in generating vibrant and visually pleasing scenes.

# Future Work

1. Reduce reliance on LLM

   - LLMs may not be able to handle numerical tasks accurately

   - Constraint solving is a highly complex task

   - Explore different shape parsers and constraint solvers

   → More realistic and physically plausible scene

# Future Work

2.  Employ object generation

    - Current curated database has only ~50K objects (~23 GB)

    - Original database has ~800K (~8 TB) → impractical to keep enlarging the database

    - With object generation model that can synthesize new objects in real time:

        - If the weighted object score is lower than a certain threshold → generate object

        - Acceptable to keep using a smaller database

        - No problem if the database does not have the target object

# Future Work

3. Continue our journey on automatic XR testing
   - Scenethesis is just the first step in the final automated pipeline
   - With prompts tailor-made for specific target environments, we can generate a diverse set of environments
   - E.g., after figuring out how to make a smartphone recognize a virtual environment as a real one, we can freely test AR applications using our generated scenes
   - Then perhaps we can use a Vision Language Model to analyze key frames and produce an evaluation report

# Q&A