



NeoBabel: A Multilingual Open Tower for Visual Generation

Mohammad Mahdi Derakhshani², Dheeraj Varghese², Marzieh Fadaee^{♦1},
and Cees G. M. Snoek^{♦2}

¹Cohere Labs, ²University of Amsterdam

Corresponding authors: m.m.derakhshani@uva.nl, marzieh@cohere.com

Abstract

Text-to-image generation advancements have been predominantly English-centric, creating barriers for non-English speakers and perpetuating digital inequities. While existing systems rely on translation pipelines, these introduce semantic drift, computational overhead, and cultural misalignment. We introduce NEOBABEL, a novel multilingual image generation framework that sets a new Pareto frontier in performance, efficiency and inclusivity, supporting six languages: *English, Chinese, Dutch, French, Hindi, and Persian*. The model is trained using a combination of large-scale multilingual pretraining and high-resolution instruction tuning. To evaluate its capabilities, we expand two English-only benchmarks to multilingual equivalents: m-GenEval and m-DPG. NEOBABEL achieves state-of-the-art multilingual performance while retaining strong English capability, scoring 0.75 on m-GenEval and 0.68 on m-DPG. Notably, it performs on par with leading models on English tasks while outperforming them by +0.11 and +0.09 on multilingual benchmarks, even though these models are built on multilingual base LLMs. This demonstrates the effectiveness of our targeted alignment training for preserving and extending cross-lingual generalization. We further introduce two new metrics to rigorously assess multilingual alignment and robustness to code-mixed prompts. Notably, NEOBABEL matches or exceeds English-only models while being 2–4× smaller. We release an open toolkit, including all code, model checkpoints, a curated dataset of 124M multilingual text-image pairs, and standardized multilingual evaluation protocols, to advance inclusive AI research. Our work demonstrates that multilingual capability is not a trade-off but a catalyst for improved robustness, efficiency, and cultural fidelity in generative AI.

	Website	https://Neo-Babel.github.io
	Code	https://github.com/mnderakhshani/NeoBabel
	Models	https://hf.co/mderakhshani/NeoBabel
	Pretraining Data	https://hf.co/datasets/mderakhshani/NeoBabel-Pretrain
	Instruction Data	https://hf.co/datasets/mderakhshani/NeoBabel-Instruct
	Evaluation Data	https://hf.co/datasets/mderakhshani/NeoBabel-Eval

[♦]Principal senior advisors.

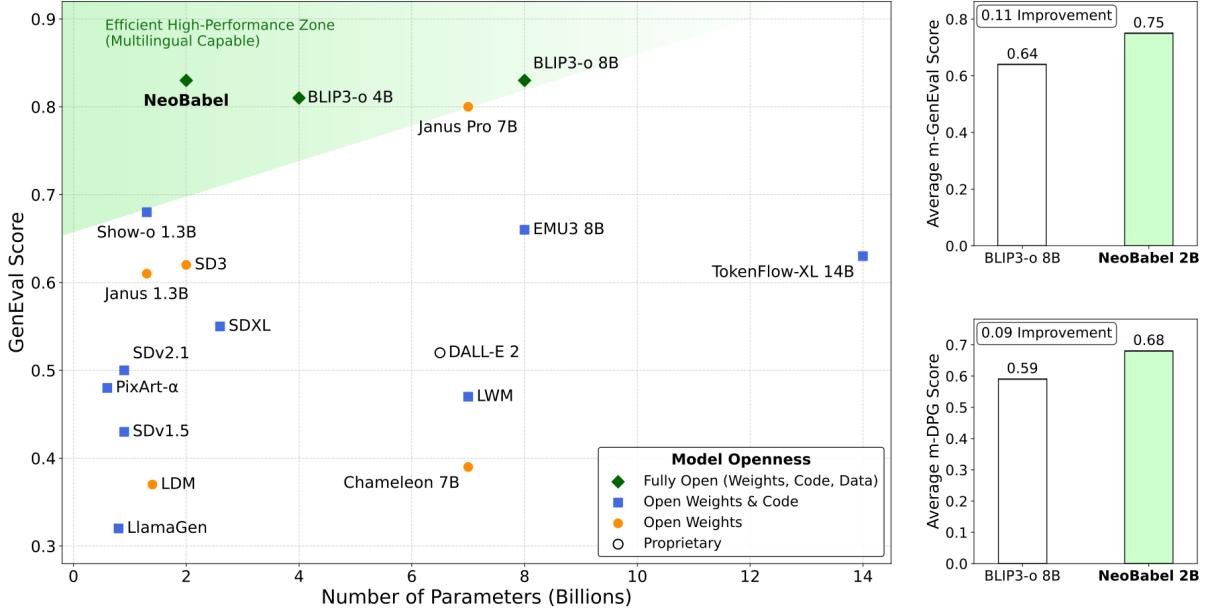


Figure 1: **NeoBabel establishes a new Pareto frontier in multilingual image generation performance, efficiency, and inclusivity.** Left: GenEval English-only scores show that NEOBABEL matches state-of-the-art models despite being 2–4× smaller. Right: On our multilingual benchmark extensions, m-GenEval and m-DPG, NEOBABEL outperforms the second-best model, demonstrating strong multilingual generalization. NEOBABEL is fully open (weights, code, data) and supports six languages with consistent cross-lingual performance.

1 Introduction

Recent advances in diffusion models and large-scale vision-language pretraining have revolutionized text-to-image generation, enabling the creation of high-quality images from natural language descriptions [Rombach et al., 2022; Peebles & Xie, 2023; Bao et al., 2023; Chen et al., 2024a; Xie et al., 2023; Wu et al., 2023a; Lipman et al., 2022; Xie et al., 2025a; Qin et al., 2025; Zhang et al., 2023; Seawead et al., 2025]. Despite these remarkable capabilities, the field suffers from a critical limitation: an overwhelming reliance on English as the primary—and often exclusive—input language [Ramesh et al., 2022; Xie et al., 2025b; Team, 2024]. This monolingual bias creates substantial barriers for the billions of users who communicate in other languages, fundamentally restricting global access to state-of-the-art generative AI technologies [Bassignana et al., 2025; Peppin et al., 2025]. The consequences of this linguistic limitation extend far beyond mere inconvenience. As text-to-image systems become integral to education, creative industries, art, and journalism, the lack of native multilingual support perpetuates existing digital divides and cultural inequities [Liu et al., 2023; Rege et al., 2025]. Non-English speakers are forced to navigate through translation layers that not only introduce friction but also risk losing the nuanced meanings and cultural contexts that make their creative expressions unique [Kannen et al., 2024; Friedrich et al., 2024]. Building truly multilingual models, like we do in this paper, is therefore not merely a technical challenge but an ethical imperative, one that ensures equitable access to generative AI while preserving linguistic diversity and cultural authenticity in the digital age.

Existing approaches to multilingual image generation typically employ a translation-first strategy, converting non-English prompts to English before processing. While this appears pragmatic, it introduces a cascade of problems that fundamentally compromise the user experience [Kreutzer et al.,

2025; Li et al., 2025b; Bafna et al., 2025]. The computational overhead of chaining translation and generation models effectively doubles inference time, creating prohibitive delays for real-time applications, thereby further disadvantaging non-English speakers. Most critically, this approach suffers from semantic drift—the systematic loss of culturally specific meanings and linguistic subtleties [Cohn-Gordon & Goodman, 2019; Vanmassenhove et al., 2019; Beinborn & Choenni, 2020]. For instance consider the Dutch term “*gezellig*” which encompasses a complex blend of coziness, conviviality, and belonging and has no direct English equivalent. When forced through translation, such rich cultural concepts are inevitably flattened or distorted, resulting in generated images that fail to capture the intended meaning. The fundamental issue lies deeper than mere translation accuracy [Wein & Schneider, 2023; Singh et al., 2024; Salazar et al., 2025].

Current vision-language architectures treat multilingual support as an afterthought, forcing diverse linguistic communities to conform to English-centric models rather than developing systems that natively understand and respect linguistic diversity. This design philosophy not only limits accessibility but also wastes the potential benefits of multilingual training, which could enhance model robustness, cross-cultural understanding, and generalization capabilities across different linguistic and cultural contexts [Ji et al., 2024; Faisal & Anastasopoulos, 2024; Dash et al., 2025; Shimabucoro et al., 2025]. These challenges demand a paradigm shift toward native multilingual understanding in text-to-image generation. The primary obstacle remains the scarcity of high-quality, culturally annotated visual-linguistic datasets for non-English languages. Even with adequate data, significant technical barriers persist: establishing robust cross-lingual concept alignment, modeling typological variations across language families, and preserving culture-specific semantics during generation. Overcoming these limitations is critical for transitioning from mere translation-based approaches to systems with genuine multilingual competence.

This paper introduces NEOBABEL, a novel multilingual image generation framework that represents the first scalable solution for direct text-to-image synthesis across six languages. Through meticulous curation of high-quality multilingual vision-language datasets and end-to-end training, NEOBABEL establishes direct cross-lingual mappings between textual descriptions and visual outputs across all supported languages. This approach not only removes translation dependencies but also maintains crucial cultural and linguistic specificity in the generated images. Our model demonstrates that multilingual capability isn’t a trade-off but rather a catalyst for improved model performance.

Our work addresses three key questions: 1) *How can we train a single model to handle multiple languages effectively?* 2) *Does multilingual training degrade performance in high-resource languages like English?* and 3) *Can a unified model outperform language-specific or translation-based approaches?* To answer these, we introduce a progressive training pipeline that combines large-scale multilingual pretraining with high-resolution instruction tuning. We evaluate NEOBABEL on m-GenEval and m-DPG, our multilingual extensions of GenEval [Ghosh et al., 2023] and DPG-Bench [Hu et al., 2024], and introduce two new metrics, Cross-Lingual Consistency (CLC) and Code Switching Similarity (CSS), to quantify multilingual performance.

As shown in Figure 1, NEOBABEL matches the performance of state-of-the-art English-only models while being 2–4× smaller. Here, we report English-only results for fair comparison, as prior work evaluates only in English. Furthermore, NEOBABEL maintains strong generation quality in all six supported languages. For instance, on the m-GenEval benchmark, it achieves a new state-of-the-art score of 0.75—an improvement of 0.11 over the very recent BLIP3-o 8B model (0.64) [Chen et al., 2025a]. Similarly, on m-DPG, it scores 0.68, outperforming BLIP3-o 8B by 0.09. These

results demonstrate that strong multilingual generation is achievable without resorting to large-scale models or sacrificing output quality.

To summarize, we make the following key contributions:

1. **A novel multilingual training framework.** We introduce a novel multilingual training framework that establishes new state-of-the-art performance in cross-lingual image generation. Our approach achieves language-agnostic understanding by directly mapping prompts from any supported language to visual concepts without requiring translation, while maintaining performance parity that matches or exceeds English-only models across all languages. This unified architecture delivers significant operational efficiency gains by eliminating the need for separate translation infrastructure, enabling single-model deployment that reduces both computational overhead and system complexity. The unified architecture delivers significant efficiency improvements, processing multilingual prompts 2.8x faster than translation-then-generation pipelines while using 59% less memory which is critical for real-world deployment scenarios. To train the unified model, we introduce a data curation pipeline that prepares multilingual image-text pairs for both pretraining and instruction tuning.
2. **Comprehensive multilingual benchmark and metrics.** We introduce the first standardized framework for evaluating multilingual image generation, addressing critical gaps in existing benchmarks. Our protocol includes: (1) extended versions of GenEval [Ghosh et al., 2023] and DPG-Bench [Hu et al., 2024], referred to as m-GenEval and m-DPG, across six languages, enabling direct comparison between native multilingual and translation-based approaches; and (2) two novel metrics—Cross-Lingual Consistency (CLC) and Code-Switching Similarity (CSS), to quantify semantic alignment and robustness to mixed-language prompts (see Figure 8). CLC measures image equivalence across languages using EVA-CLIP [Sun et al., 2023b] and DINOv2 [Oquab et al., 2023] embeddings, while CSS evaluates real-world code-switching scenarios. NEOBABEL achieves state-of-the-art multilingual performance while maintaining strong English capabilities. Notably, it matches the English results of leading multilingual models while outperforming them by +0.11 and +0.09 on multilingual benchmarks—despite those models being built on multilingual base LLMs. This positions NEOBABEL as a strong foundation for future research in equitable, culturally adaptive generative AI.
3. **Open toolkit for inclusive research.** We release a comprehensive research toolkit comprising NEOBABEL model checkpoints trained on six languages (English, Chinese, Dutch, French, Hindi, and Persian), a systematically curated dataset of 124M multilingual text-image pairs with quality-controlled translations, and a complete reproducibility package including training scripts, hyperparameter configurations, and standardized evaluation protocols. Our framework is designed to be easily extensible to additional languages, thanks to a scalable training pipeline, with validation metrics and benchmarking guidelines that support systematic comparison of multilingual generation across research groups.

In the following sections, we present the details of NEOBABEL, including its architecture (Section 2), multilingual datasets (Section 3), progressive training stages (Section 4), and multilingual evaluation suite (Section 5). We then provide both quantitative and qualitative evaluations (Section 6), followed by ablation studies and analysis (Section 7).

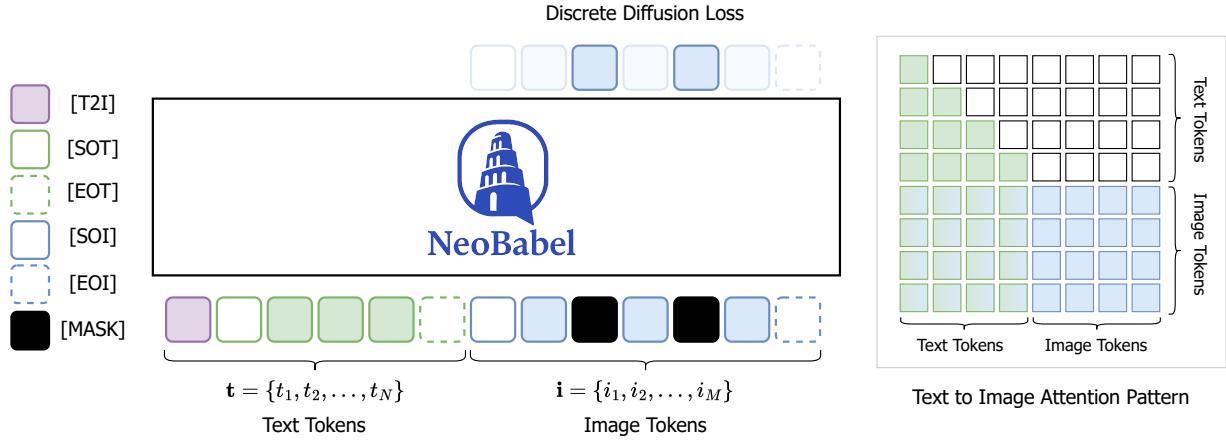


Figure 2: **NeoBabel: A Multilingual Open Tower for Visual Generation.** Regardless of modality, all input data is first tokenized and embedded into a unified input sequence. NEOBABEL then applies causal attention to text tokens and full attention within a discrete denoising diffusion framework for image tokens, ultimately generating the desired image. This design enables NEOBABEL to support a wide range of tasks, including text-to-image generation, text-guided inpainting and extrapolation, as well as cross-lingual image generation.

2 NeoBabel Architecture

We first outline the core architectural components of NEOBABEL, including its multilingual transformer backbone (Section 2.1), training objectives (Section 2.2), and the multilingual model merging strategy (Section 2.3) designed to enhance generation quality across diverse linguistic settings.

2.1 Model Architecture

Our architecture’s core components, a multilingual tokenizer and transformer backbone, are specifically optimized for efficient, scalable cross-lingual image generation, supporting seamless processing across diverse languages and image types. Figure 2 provides an overview of the NEOBABEL architecture.

2.1.1 Tokenizers

Text Tokenization. For textual input, we adopt the tokenizer of the pretrained multilingual large language model Gemma-2 [Gemma Team et al., 2024] without any modifications. This approach maintains compatibility with multilingual inputs while utilizing proven tokenization methods from language modeling.

Image Tokenization. For image input, we leverage the MAGVIT-v2 quantizer [Yu et al., 2023] retrained by Show-o [Xie et al., 2025b] on 25 million images. This lookup-free quantizer learns a discrete codebook of size $K=8,192$ and encodes 256×256 resolution images into 16×16 grids of discrete tokens. The quantization approach supports efficient downstream training and generation while preserving fine-grained visual details.

2.1.2 Transformer Backbone

As we build upon the pretrained multilingual large language model (LLM) Gemma-2 [Gemma Team et al., 2024], we maintain its overall transformer architecture, while introducing two key modifications: (1) integration of a unified multimodal embedding space, and (2) modality-aware attention patterns for flexible generation. Additionally, we apply qk-norm [Henry et al., 2020] to each attention layer to enhance training stability and convergence.

Unified Multimodal Embedding and Prompt Design. To enable seamless multimodal learning, we extend the LLM’s embedding table with 8,192 new learnable embeddings for discrete image tokens, allowing the model to process image inputs natively without architectural changes. Both text and image tokens are embedded in a shared space, enabling the model to learn cross-modal compositionality and semantic alignment. We represent all tasks including text-to-image generation as unified autoregressive sequences. Given a tokenized image-text pair, text and image tokens are concatenated into a single sequence. Special tokens such as [T2I], [SOT], [EOT], [SOI], and [EOI] explicitly mark task type and modality boundaries, enabling the model to disambiguate different modalities and tasks through prompting alone. This design simplifies the training pipeline by removing the need for modality-specific components or task-specific heads, allowing for flexible, scalable, and unified multimodal generation.

Modality-Aware Attention Patterns. To accommodate the differing structural needs of text and image modalities, we employ a hybrid attention mechanism. Text tokens are modeled with causal attention to preserve autoregressive language modeling capabilities. Image tokens, in contrast, are modeled using full bidirectional attention, allowing rich interactions that are critical for high-fidelity image synthesis. When both modalities are present, attention masks are dynamically configured so that image tokens can fully attend to text tokens and preceding image tokens, enabling coherent, contextually grounded generation.

2.2 Training Objective

The model is trained on sequences composed of both textual and visual tokens, where text tokens act as a prefix and visual tokens form the postfix. We do not apply any learning objective to the text tokens; the loss is computed solely over the visual tokens.

Let $\mathbf{t} = \{t_1, t_2, \dots, t_N\}$ denote the text tokens and $\mathbf{i} = \{i_1, i_2, \dots, i_M\}$ denote the image tokens, forming a full input sequence $[\mathbf{t}; \mathbf{i}]$. During training, we randomly select a subset $\mathcal{J} \subset \{1, \dots, M\}$ of image token indices to be masked. The corresponding masked sequence is denoted by \mathbf{i}_* , where i_j is replaced with a special [MASK] token for all $j \in \mathcal{J}$. The model is trained to reconstruct the original visual tokens at the masked positions by conditioning on the full input sequence of text tokens and (partially masked) image tokens. The objective is defined as:

$$\mathcal{L} = \sum_{j \in \mathcal{J}} \log p_\theta(i_j | \mathbf{t}, \mathbf{i}_*), \quad (1)$$

where $p_\theta(\cdot)$ is the model’s predicted distribution over image codebook entries, parameterized by θ . The loss is only applied to the masked image tokens in \mathcal{J} . We follow the masking strategy introduced by Xie et al. [2025b], randomly masking a fixed ratio of visual tokens within each training sample. To further improve generation controllability, we incorporate classifier-free guidance [Ho & Salimans, 2022] by replacing the conditioning text with a null string with some probability during training.

2.3 Multilingual Model Merging

To enhance generalization and stability of multilingual image generation models, we adopt model merging techniques that combine multiple checkpoints from the training trajectory. Let $\{M_i\}_{i=1}^N$ denote a sequence of N model checkpoints and $\{w_i\}_{i=1}^N$ their corresponding non-negative weights. The merged model \widehat{M} is defined as a convex combination:

$$\widehat{M} = \sum_{i=1}^N \alpha_i M_i \quad \text{where} \quad \alpha_i = \frac{w_i}{\sum_{j=1}^N w_j}. \quad (2)$$

This formulation allows the merged model to interpolate within the solution space spanned by the selected checkpoints, potentially improving generalization on unseen prompts and enhancing robustness to overfitting. We consider three widely used weighting strategies for this purpose, each reflecting different assumptions about model evolution during training. The comparative results and analysis of these approaches are presented ablation studies section.

Simple Moving Average (SMA) assigns equal weight to all checkpoints. It is defined as:

$$M_{\text{avg}} = \frac{1}{N} \sum_{i=1}^N M_i. \quad (3)$$

SMA is simple, stable, and particularly effective when applied in the later stages of training where model weights exhibit minimal drift. Prior work [Li et al., 2025c] found SMA to perform robustly due to this stabilization.

Exponential Moving Average (EMA) emphasizes recent checkpoints by applying exponentially decaying weights. It is computed recursively as:

$$M_{\text{avg}}^{(i)} = \alpha M_i + (1 - \alpha) M_{\text{avg}}^{(i-1)}, \quad i \in [2, N]. \quad (4)$$

The decay factor $\alpha \in (0, 1)$ controls the trade-off between recency and stability. EMA adapts more quickly to recent model dynamics but is sensitive to noise if weights are unstable.

Weighted Moving Average (WMA) assigns custom, possibly increasing weights to later checkpoints. The merged model is computed using the normalized form:

$$M_{\text{avg}} = \sum_{i=1}^N \frac{w_i}{w_{\text{sum}}} M_i, \quad \text{where} \quad w_{\text{sum}} = \sum_{i=1}^N w_i. \quad (5)$$

This general formulation allows flexibility in how much importance is placed on each checkpoint. In our case, we use $w_i = i$ to emphasize later-stage models.

3 NeoBabel Multilingual Datasets

3.1 Data Curation Pipeline

Multilingual multimodal data remains scarce, especially compared to the abundance of English-centric resources. This imbalance poses a significant barrier to training and evaluating models that

Original English-Only Dataset				NEOBABEL Multilingual Expansion		
Dataset	Image Source	Caption Source	Size	Recapitulation	Translation	New Size
ImageNet 1K	Web	Class labels	1M	—	✓	6M
CC12M	Web	Alt-text (noisy)	12M	✓	—	12M
SA-1B	Photography	LLaVA	10M	✓	—	10M
LAION-Aesthetic	Web	Alt-text (noisy)	12M	✓	✓	72M
JourneyDB	Synthetic	GPT-3.5	4M	✓	✓	24M
BLIP3-o Instruct	Web + Synthetic	GPT-4o / human	60K	—	✓	360K
				39M		
				124M		

Table 1: **NeoBabel multilingual datasets**, detailing their English-only data source, image origin, caption format, and size. Our multilingual expansion covers model-generated recaptioning, translation into multiple languages, or both. Our expansions increase the total size from 39M to 124M image–caption/label pairs. In the remainder of this paper, all modified datasets are prefixed with m- to denote their expanded form.

can understand grounded language across diverse linguistic contexts. To address this gap, we curate and augment several multilingual datasets by translating and recaptioning existing image-caption pairs into six target languages: English, Chinese¹, Dutch, French, Hindi, and Persian. We summarize the datasets curated in Table 1. At the core of our approach is a multilingual captioning pipeline designed to ensure both semantic richness and linguistic diversity. We begin by generating a detailed English caption for each image using InternVL [Chen et al., 2024c], prompted with a simple instruction: “Describe this image in detail in English.” This step guarantees comprehensive coverage of the visual content.

To preserve quality and consistency across languages, we implement a multi-step post-processing and filtering stage based on four strategies:

- **Length filtering:** Remove captions that are too short (e.g., fewer than 5 tokens) or excessively long (e.g., more than 500 tokens).
- **Language validation:** Detect and discard captions containing non-English phrases or corrupted outputs using language identification tools. We use the fastText language identification model trained on 176 languages [Joulin et al., 2016]. We discard any caption not classified as English with a confidence score above 90%.
- **Visual-text mismatch filtering:** Discard captions that do not align with visual content, measured via auxiliary vision-language models (e.g., using VQAScore). Specifically, we leverage MolMo-72B [Deitke et al., 2025] deployed with vLLM [Kwon et al., 2023], formulating the task as a binary structured prediction (yes/no) via vLLM’s output interface.
- **Toxicity and NSFW filtering:** Discard samples using the LAION-5B NSFW classifier [Schuhmann et al., 2022] to ensure safe visual content before captioning, assuming high likelihood of appropriateness in the resulting captions.

Once high-quality English captions are obtained, we translate them into five target languages using the NLLB model [Costa-Jussà et al., 2022] for the pretraining datasets, and the Gemini Experimental model (gemini-2.0-flash-lite) for the instruction tuning datasets. This separation ensures

¹Throughout this work ‘Chinese’ refers to Simplified Chinese.

high translation coverage at scale during pretraining, while leveraging higher-quality outputs for instruction-tuned data. Using English as a pivot allows us to take advantage of strong captioning performance in high-resource settings while ensuring consistent semantic content across all languages. This approach not only amplifies the linguistic diversity of our dataset but also maintains alignment between captions, which is critical for multilingual training and evaluation. Ultimately, this step plays a central role in constructing a high-quality, language-balanced multimodal resource—an essential step toward more inclusive and globally-relevant vision-language models.

3.2 NeoBabel Pretraining Data

The previous section described the overall pipeline and transformation steps, and next we detail the data sources and multilingual adaptations used to train the model. We use a diverse collection of image-text datasets to build strong multilingual visual-language alignment combining real-world and synthetic image sources. While the images are drawn from established, high-quality datasets, the accompanying captions have been significantly enriched through our recaptioning and multilingual translation pipeline—resulting in a more diverse, detailed, and valuable resource for future multilingual generative models.

m-ImageNet-1K: The original English class labels are translated into five more languages to obtain a total of six target languages, forming multilingual textual prompts for class-conditional image generation.

m-SA-1B and m-CC12M: We incorporate 22 million image-caption pairs in English from SA-1B [Kirillov et al., 2023] and CC12M [Changpinyo et al., 2021]. These datasets provide rich natural image-caption pairs and enhance visual diversity. The texts are enhanced through our recaptioning pipeline described in Section 3.1.

m-LAION-Aesthetic: A subset of the LAION dataset including 12M image-text pairs² is enhanced and translated, yielding approximately 72 million image-caption pairs for a total of six languages.

m-JourneyDB: This synthetic dataset consists of 4 million high-quality images generated by the Midjourney model [Sun et al., 2023a]. We apply the same recaptioning and translation pipeline to generate 24 million image-caption pairs for our six languages.

Combining all sources, the final pretraining dataset contains approximately 124 million image-text pairs across six languages, covering diverse domains and visual aesthetics.

3.3 NeoBabel Instruction Tuning Data

Here we describe our datasets and mixing strategies used for instruction tuning. This phase reuses two datasets introduced earlier and adds a smaller but higher-quality dataset focused on multimodal instruction tuning:

m-LAION-Aesthetic and m-JourneyDB: Our setup continues to use the LAION-Aesthetic and JourneyDB datasets, as extended in the pretraining data.

²<https://huggingface.co/datasets/dclure/laion-aesthetics-12m-umap>

m-BLIP3o-Instruct: An instruction-focused dataset introduced by Chen et al. [2025a], containing multimodal instruction samples, also translated into six languages for multilingual supervision.

All images are resized to 512×512 . While the images are drawn from established, high-quality sources, most accompanying texts have been significantly enriched or rewritten, resulting in a more valuable and linguistically diverse dataset for instruction tuning and multilingual generation.

4 NeoBabel Training Stages: Learning Progression

NEOBABEL is trained using a staged learning framework consisting of three progressive pretraining stages (Section 4.1) followed by two instruction tuning stages (Section 4.2).

4.1 Progressive Pretraining

Our pretraining includes three stages, progressively scaling from basic visual understanding to advanced multilingual image generation:

Stage 1 – Pixel Dependency Learning: The model initially learns foundational visual representations using m-ImageNet-1K. Class-conditional image generation is guided by translated class labels, enabling the model to form robust image token embeddings and capture pixel-level dependencies for high-fidelity output.

Stage 2 – Scaling Alignment with Large-Scale Multilingual Data: Using weights from the first stage, the model is fine-tuned on 22 million English-only image-caption pairs (from m-SA-1B and m-CC12M) and 72 million translated samples from m-LAION-Aesthetic. This stage strengthens the model’s grounding in natural image-text alignment while developing multilingual capabilities through broad cross-lingual exposure.

Stage 3 – Refined Multilingual Pretraining: In the final stage, the model is trained on 96 million multilingual image-text pairs derived from m-LAION-Aesthetic and m-JourneyDB. The training balances high-quality real-world aesthetic data with diverse, synthetic images to improve generalization across languages, domains, and modalities.

4.2 Progressive Instruction Tuning

Following pretraining, the model advances to instruction tuning, where the focus shifts from unsupervised representation learning to explicit task-guided adaptation, refining its ability to interpret and execute complex, multilingual instructions through our curated datasets and progressive exposure to prompt-driven generation in two stages:

Stage 1 – Initial Multilingual Instruction Alignment: To build robust multilingual instruction-following capabilities at high resolution, the model is first trained with a diverse mixture of the three datasets described above. In this stage, training samples are drawn from m-LAION-Aesthetic, m-JourneyDB, and m-BLIP3o-Instruct using mixing weights α_1 , α_2 , and α_3 , respectively, such that $\alpha_1 + \alpha_2 + \alpha_3 = 100$. A higher α_1 and moderate α_2 prioritize real-world and aesthetic content, while a smaller α_3 introduces early exposure to instruction-rich samples. This balance helps the model

learn cross-lingual, cross-modal grounding without overwhelming it with complex prompts in the early stages.

Stage 2 – Instruction Refinement: In the second stage, we adjust the mixing weights to emphasize instruction-rich and synthetic supervision. Specifically, α_2 and α_3 are increased to draw more heavily from m-JourneyDB and m-BLIP3o-Instruct, while α_1 is decreased to reduce reliance on LAION-based content. This curriculum-style shift enables the model to refine its instruction-following capabilities using complex multilingual prompts and high-quality synthetic images. The increased semantic richness improves the model’s generalization to both benchmark instruction tasks and open-ended generation scenarios.

Each stage is trained for 500k steps (except the final stage of instruction tuning with 200k) using the AdamW optimizer and cosine learning rate decay. The learning rate is set to $1e-4$ during pretraining and adjusted during instruction tuning. We gradually increase prompt sequence length and resolution from 128 to 512 and from 256×256 to 512×512 respectively. The vocabulary and codebook sizes are fixed across all stages. Full hyperparameter settings for each pretraining and instruction tuning stage are summarized in the Appendix.

5 Multilingual Evaluation of Image Generation

Existing image generation benchmarks are mostly English-centric, failing to capture cross-lingual performance. To resolve this limitation, we introduce a multilingual evaluation suite that extends established (English-only) benchmarks to cover six diverse languages and introduces new evaluation metrics for assessing cross-lingual visual consistency. This section outlines our multilingual evaluation suite (Section 5.1) and multilingual evaluation metrics (Section 5.2).

5.1 Multilingual Evaluation Suite

We assess the image generation capabilities of NEOBABEL using two complementary benchmarks: GenEval [Ghosh et al., 2023] and DPG-Bench [Hu et al., 2024]. GenEval offers a structured evaluation of prompt-to-image alignment across six compositional dimensions: *single object*, *two objects*, *counting*, *colors*, *position*, and *color attribute*. In contrast, DPG-Bench targets general-purpose generation with open-ended, diverse prompts that test broader semantic understanding. However, both benchmarks are English-only and fail to capture multilingual generative performance.

As part of our multilingual evaluation suite, we introduce **m-GenEval** and **m-DPG**, multilingual extensions of the original benchmarks. All prompts are translated into five additional languages: **Chinese**, **Dutch**, **French**, **Hindi**, and **Persian**, using the Gemini Experimental model, followed by human verification and manual corrections to ensure semantic fidelity and linguistic fluency. Together with the paper, we publicly release m-GenEval and m-DPG to promote inclusive and realistic evaluation of multilingual text-to-image models and support broader community adoption.

5.2 Multilingual Evaluation Metrics

To complement the multilingual benchmarks introduced above, we introduce two metrics that assess how well generative models preserve visual and semantic consistency across languages. Existing

evaluations focus on monolingual alignment, overlooking whether models produce consistent outputs across languages or under mixed-language inputs. To address this, we introduce two scores to assess cross-lingual consistency and robustness under intra-prompt language mixing. Together, these metrics provide a more diagnostic view of multilingual generation performance.

Cross-Linguistic Consistency (CLC). To evaluate whether multilingual models generate semantically consistent and faithful outputs across languages, we introduce the CLC score. Multilingual image generation models should produce visually similar outputs when given semantically equivalent prompts, regardless of the input language. Measuring this consistency is crucial for understanding how well the model aligns its multilingual text inputs with the corresponding visual outputs, which reflects the quality of its cross-lingual grounding. We evaluate in a multilingual setting consisting of P prompts, each paired with L language variations, forming a parallel dataset. For each prompt $p \in \{p_i\}_{i=1}^P$, we generate K images (one per language), resulting in $L \times K$ images. Let x_i denote an image and $f(x_i) \in \mathbb{R}^d$ its corresponding embedding obtained from a vision encoder.

To measure consistency, we treat the K images generated from the English version of the prompt as the reference set \mathcal{R}_p , and the remaining $(L - 1) \times K$ images generated from other languages as the target set \mathcal{T}_p . The core idea is that if the model is truly language-agnostic in its understanding, images generated from non-English prompts should be visually similar to those generated from the English prompt. The CLC Score for prompt p is computed by averaging the cosine similarity between all reference and non-reference embeddings:

$$\text{CLC}_p = \frac{1}{|\mathcal{R}_p| \cdot |\mathcal{T}_p|} \sum_{x_i \in \mathcal{R}_p} \sum_{x_j \in \mathcal{T}_p} \cos(f(x_i), f(x_j)). \quad (6)$$

Finally, the overall CLC score is obtained by averaging CLC_p over all prompts P . For evaluation, we use m-DPG prompts and compute embeddings with two strong vision encoders, EVA-CLIP [Sun et al., 2023b] and DINOv2 [Oquab et al., 2023], to ensure robustness across different feature representations. This metric provides a quantitative measure of how well multilingual generation models maintain semantic and visual alignment across languages.

Code-Switching Similarity (CSS). Real-world multilingual communication frequently involves code switching, i.e., interleaving of multiple languages within a single utterance. Therefore, a well-aligned multilingual model should demonstrate robustness not only to monolingual prompts but also to mixed-language inputs, capturing the inherent complexity and variability of natural language. Code switching often increases perplexity and degrades performance in language models; however, its impact on image generation remains largely unexplored. To evaluate this, we introduce the CSS Score, which quantifies visual consistency under intra-prompt language variation. Given a set of reference prompts composed entirely in English, we construct two variants per prompt for each of the $L-1$ non-English target languages: (1) **English-First (EF)**: the first half of the prompt remains in English while the second half is translated into the target language, and (2) **English-Second (ES)**: the first half is translated while the second half remains in English.

For each prompt $p \in \{p_i\}_{i=1}^P$, we generate a single reference image x_{ref} from the original English prompt and $L - 1$ code-switched images: $x_{\text{EF}}^{(l)}$ and $x_{\text{ES}}^{(l)}$ for each target language l . Each image is encoded into an embedding $f(x) \in \mathbb{R}^d$ using a vision encoder. The Code Switching Similarity (CSS) score for each prompt is computed by measuring the average cosine similarity between the reference

Method	⌚	Type	Params.	Single Object	Two Object	Counting	Colors	Position	Color Attribute	Overall
LlamaGen	×	G	0.8B	0.71	0.34	0.21	0.58	0.07	0.04	0.32
LDM	×	G	1.4B	0.92	0.29	0.23	0.70	0.02	0.05	0.37
SDv1.5	×	G	0.9B	0.97	0.38	0.35	0.76	0.04	0.06	0.43
PixArt-alpha	×	G	0.6B	0.98	0.50	0.44	0.80	0.08	0.07	0.48
SDv2.1	×	G	0.9B	0.98	0.51	0.44	0.85	0.07	0.17	0.50
DALL-E 2	×	G	6.5B	0.98	0.66	0.49	0.77	0.10	0.19	0.52
SDXL	×	G	2.6B	0.98	0.74	0.39	0.85	0.15	0.23	0.55
SD3	×	G	2B	0.98	0.74	0.63	0.67	0.34	0.36	0.62
CoDI	×	U&G	-	0.89	0.16	0.16	0.65	0.02	0.01	0.31
Chameleon	×	U&G	7B	-	-	-	-	-	-	0.39
LWM	○	U&G	7B	0.93	0.41	0.46	0.79	0.09	0.15	0.47
SEED-X	○	U&G	17B	0.97	0.58	0.26	0.80	0.19	0.14	0.49
Janus	○	U&G	1.3B	-	-	-	-	-	-	0.61
TokenFlow	○	U&G	14B	-	-	-	-	-	-	0.63
EMU3	○	U&G	8B	-	-	-	-	-	-	0.66
Show-o	×	U&G	1.3B	0.98	0.80	0.66	0.84	0.31	0.50	0.68
Janus-Pro	○	U&G	7B	-	-	-	-	-	-	0.80
BLIP3-o	○	U&G	4B	-	-	-	-	-	-	0.81
BLIP3-o	○	U&G	8B	-	-	-	-	-	-	0.83
NEOBABEL	✓	G	2B	1.00	0.91	0.62	0.91	0.81	0.77	0.83

Table 2: **English-only GenEval benchmark comparison.** NEOBABEL achieves the highest overall score, outperforming larger models on tasks requiring compositional reasoning and fine-grained prompt-image alignment. Symbol legend: ⌚ denotes multilingual generation capability, with ✓ indicates a full multilingual capability, ○ represents partial multilingual capability (i.e. bilingual or multilingual to a limited extent), and × denotes monolingual models.

embedding $f(x_{\text{ref}})$ and the embeddings from the EF and ES variants:

$$\text{CSS}^{\text{EF}} = \frac{1}{L-1} \sum_{l=1}^{L-1} \cos(f(x_{\text{ref}}), f(x_{\text{EF}}^{(l)})), \quad \text{CSS}^{\text{ES}} = \frac{1}{L-1} \sum_{l=1}^{L-1} \cos(f(x_{\text{ref}}), f(x_{\text{ES}}^{(l)})). \quad (7)$$

The final CSS scores are obtained by averaging across all prompts:

$$\text{CSS}^{\text{EF}} = \frac{1}{P} \sum_{p=1}^P \text{CSS}_p^{\text{EF}}, \quad \text{CSS}^{\text{ES}} = \frac{1}{P} \sum_{p=1}^P \text{CSS}_p^{\text{ES}}. \quad (8)$$

To assess how well models preserve semantic consistency under intra-prompt code switching, we report both CSS^{EF} and CSS^{ES} , using embeddings from EVA-CLIP [Sun et al., 2023b] and DINOv2 [Oquab et al., 2023] computed on m-DPG prompts.

6 Results and Discussions

We evaluate NEOBABEL on our multilingual extension of standard benchmarks, including m-GenEval and m-DPG, to assess performance across languages both quantitatively and qualitatively.

Baselines. We evaluate our model against a diverse range of baselines, which we group into two categories: generative-only models (G) and unified models (U&G). The generative models are

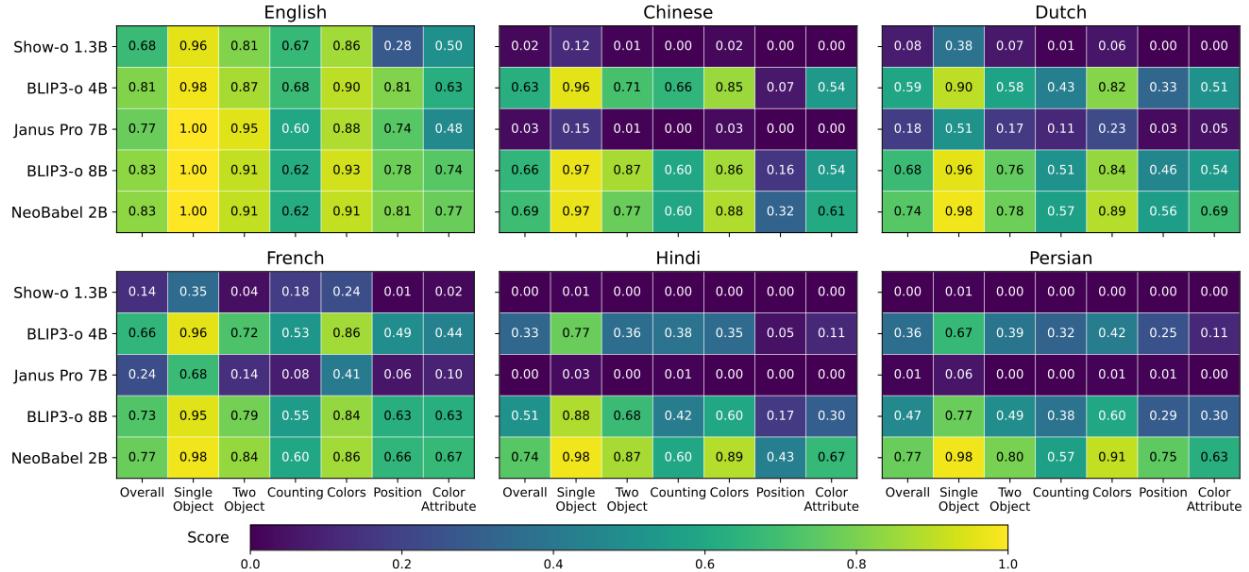


Figure 3: m-GenEval benchmark comparison. Models such as Janus Pro and BLIP3-o rely on multilingual base LLMs but are trained solely on English image-generation data, leading to a sharp performance drop in non-English languages. In contrast, NEOBABEL maintains strong and consistent results across all six languages, demonstrating robust cross-lingual generalization. Here baseline models are ordered by parameter count.

designed exclusively for text-to-image generation, without any visual understanding components. This category includes LlamaGen [Sun et al., 2024], LDM [Rombach et al., 2022], SDv1.5 and SDv2.1 [Rombach et al., 2022], SDXL [Podell et al., 2023], SD3 [Esser et al., 2024], DALL-E 2 [Ramesh et al., 2022], and PixArt- α [Chen et al., 2024a] models primarily optimized for high-quality and compositional image generation. In contrast, the unified models support both image generation and image understanding tasks such as captioning and visual question answering. This group includes CoDI [Tang et al., 2023], LWM [Liu et al., 2024a], SEED-X [Ge et al., 2024], Chameleon [Team, 2024], TokenFlow [Qu et al., 2025], EMU3 [Wang et al., 2024], Janus [Wu et al., 2025], Janus-Pro [Chen et al., 2025b], and BLIP3-o [Chen et al., 2025a]. Our comparison includes both small-scale and large-scale models, spanning from under 1B to over 17B parameters.

6.1 Multilingual Image Generation Performance

m-GenEval Comparison. We begin by evaluating NEOBABEL on the English prompts of the m-GenEval benchmark, with results reported in Table 2. The comparison includes both generative models (G), which focus solely on text-to-image generation, and unified models (U&G), which also support image understanding tasks such as captioning and visual question answering. Despite having only 2B parameters, NEOBABEL outperforms or matches best-performing unified models such as Janus-Pro 7B (0.77) and BLIP3-o 8B (0.83), which are significantly larger in terms of parameters. It also surpasses SD3 2B (0.62), a leading model in the generative category, achieving the highest overall score of 0.83. This performance reflects strong fine-grained and compositional prompt-image alignment particularly in challenging subcategories like color attributes and positional grounding.

In Figure 3, we further evaluate NEOBABEL across five more languages including Chinese, Dutch, French, Hindi, and Persian to assess its multilingual generalization capabilities beyond English. As can be seen, the performance gap between NEOBABEL and the strongest baselines is small

Model	Params.	English	Chinese	Dutch	French	Hindi	Persian	Overall
Show-o	1.3B	0.67	0.10	0.22	0.32	0.04	0.04	0.23
EMU3	8B	0.80	—	—	—	—	—	—
TokenFlow-XL	14B	0.73	—	—	—	—	—	—
Janus	1.3B	0.79	0.56	0.42	0.53	0.17	0.13	0.43
Janus Pro	7B	0.84	0.50	0.61	0.68	0.12	0.12	0.47
BLIP3-o	4B	0.79	0.60	0.58	0.59	0.47	0.49	0.58
BLIP3-o	8B	0.80	0.56	0.59	0.61	0.50	0.53	0.59
NEOBABEL	2B	0.75	0.70	0.69	0.70	0.63	0.65	0.68

Table 3: **m-DPG benchmark comparison.** Despite its small parameter count, NEOBABEL achieves competitive results in English and consistently outperforms all baselines across five non-English languages, demonstrating strong cross-lingual prompt understanding and image generation.

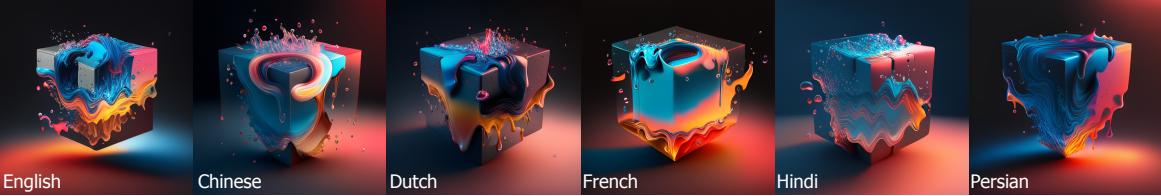
in Chinese (by 0.03). We attribute this to two factors: (i) the use of bilingual English-Chinese instruction-tuning data in models like Janus Pro, whose training setup is not publicly disclosed, and (ii) architectural choices such as BLIP3-o’s use of a frozen LLM backbone with prompt learning instead of full model adaptation. In medium-resource languages like Dutch and French, the gap widens (0.06 and 0.04 respectively), and in low-resource languages such as Hindi and Persian, NEOBABEL significantly outperforms all baselines by a large margin (up to 0.3 improvement), despite having 4× fewer parameters than Janus Pro 7B and BLIP3-o 8B. These results underscore the cross-lingual robustness and data efficiency of our multilingual instruction-tuning strategy.

m-DPG Comparison. Compared to m-GenEval, which emphasizes fine-grained attributes and atomic compositional reasoning, m-DPG focuses on a model’s ability to follow natural, descriptive multilingual prompts. It tests whether the generated images are semantically accurate, detailed, and coherent, making it a stronger indicator of real-world prompt-image alignment performance. We evaluate NEOBABEL on m-DPG to assess prompt-image alignment across 6 languages in Table 3. NEOBABEL achieves comparable performance in English (0.75), even though it uses only 2B parameters, which is far fewer than BLIP3-o (4B and 8B) and Janus Pro (7B). More importantly, NEOBABEL outperforms all baselines in the non-English settings. Existing models show notable performance drops in several languages, especially in low-resource settings (Hindi and Persian) where models such as Janus, Janus Pro, and Show-o perform poorly. In contrast to the best-performing baseline (BLIP3-o 8B), NEOBABEL consistently achieves the highest scores across all six target languages. As in m-GenEval, we observe a similar trend in m-DPG, where the performance gap widens in medium-resource languages by 0.10 in Dutch and 0.09 in French and becomes even larger in low-resource settings, with gaps of 0.13 in Hindi and 0.12 in Persian.

6.2 Qualitative Evaluation

To complement the quantitative findings, we present qualitative results from NEOBABEL across diverse prompt categories, including compositional scenes, abstract concepts, and multilingual instructions, in Figures 4 and 5 (main paper) and Figure 12 (appendix). The results show that NEOBABEL consistently generates semantically aligned and visually coherent images. Objects, layouts, and attributes are preserved across languages, demonstrating the model’s strong multilingual alignment and consistency in representing concepts.

English: The image depicts a dynamic and vibrant scene featuring an abstract, three-dimensional cube suspended against a dark background with subtle red hues at the bottom edge. The cube is composed of fluid-like material that appears to be melting into various colors such as blue, pink, orange, yellow, and hints of purple. The liquid forms swirls around its edges like waves crashing over it, creating a sense of motion and energy within the composition. Bubbles are scattered throughout both inside and outside the cube's surface, adding texture and depth to the visual effect. The top part of the cube has more intense shades of blue while transitioning smoothly towards warmer tones on one side, giving off a fiery appearance where bright pinks and yellows dominate. This contrast creates a striking interplay between cool and warm colors across different sections of the object. Light reflections highlight some parts making them appear glossy and wet, enhancing the illusion of movement through light refraction effects visible along the curves formed by the flowing substance.



Chinese: 这张图描绘了一只迷人的动画猫角色，拥有大大的绿色眼睛和富有表现力的面孔，展现出自信与好奇心。猫咪戴着一顶细致的棕色帽子，帽侧装饰有螺丝和别针等配件。它身穿精致的服装，包括一件时尚的马甲，里面似乎是一件浅色衬衫。其毛发图案在耳朵和后腿周围有橙色斑点，赋予它独特的外观。背景为白色，突出了明亮复杂的颜色、衣服的细节和角色面部的表现。尾巴轻轻地在身后弯曲。整体上，这个设计很好地传达了大胆和好奇的感觉。



Dutch: De afbeelding toont een prachtig vormgegeven glazen fles gevuld met een levendige groene vloeistof, die magisch en betoverd oogt. De fles is aan de onderkant ingewikkeld vormgegeven met sierlijke patronen die een antieke stijl suggereren. Vanuit het binnenin ontspuiten weelderige wijnranken, versierd met bladeren in verschillende vormen en maten, waarvan sommige lijken op klimopachtige bladeren die om de hals verstengd zijn. Binnenin de fles wervelt de inhoud met gloeiende, smaragdgroene energie, wat een dynamische beweging creëert alsof de fles leeft. Kleine belletjes zweven omhoog door de vloeistof en voegen diepte en dynamiek toe aan de scène. Een paar kleine voorwerpen zoals munten en stukjes papier liggen verspreid op het oppervlak onder de fles, wat een gevoel van mysterie suggereren over hun oorsprong of doel. De belichting werpt zachte schaduwen over deze elementen, wat de mystieke sfeer versterkt en tegelijkertijd de ingewikkelde details zowel binnen als buiten de fles benadrukt. Over het geheel genomen is er een fantasie-element dat doet denken aan alchemie of het maken van toverdranken, wat nieuwsgierigheid en verwondering oproept.



French: La photo montre une jeune femme aux cheveux foncés, coiffée avec élégance, ornée d'une tiare complexe qui scintille subtilement sur son front et ses tempes. Son maquillage est discret mais raffiné, mettant en valeur ses grands yeux grâce à des sourcils bien dessinés. Elle porte une robe luxueuse en tissu riche, avec des broderies élaborées et des perles sur le corsage, formant des motifs gracieux au niveau de la poitrine. Les manches ajustées de la robe sont décorées de détails floraux délicats sur les épaules, ajoutant à son élégance. Un élément remarquable comprend des décorations semblables à des plumes s'étendant de son encolure jusqu'à sa clavicule, lui donnant une allure royale, comme si elles faisaient partie intégrante du vêtement plutôt que d'être des accessoires. Son expression est dramatique et gracieuse, sur un fond sombre et serein qui met en valeur la douceur des textures.



Figure 4: **Qualitative evaluation of NEOBABEL.** Each row is based on a single concept expressed in six different languages. For clarity, we show only one of the prompts (in one language) and present six images generated from its translated prompts in the other five languages. Across all languages, NEOBABEL delivers semantically accurate and visually cohesive outputs with reliable consistency.

6.2.1 Multilingual Image Inpainting and Extrapolation

NEOBABEL enables new collaborative applications, such as a multilingual visual canvas where users can contribute prompts in their native languages to co-create coherent and expressive visual scenes.

Hindi: चित्र एक शैलीबद्ध, कार्टन शैली के कंकाल चरित्र को चित्रित करता है, जो एक जादूगर या जादूगर के रूप में पहना हुआ है, जो दृश्य के कुछ हिस्सों को उजागर करने वाले नीले प्रकाश के साथ नीले प्रकाश के उच्चारण के साथ एक अधिरे पुष्टपूर्णि के खिलाफ खड़ा है। यह आकड़ा एक ओवरसाइज हुई ड्रेस पहनता है जो अपने सिर और कंपी को पूरी तरह से कवर करता है, केवल खोपड़ी को दोनों दो बड़े आंखों के साथम से दिखाई देता है, जहां आंखें स्थित हैं। वेहरा कंकाल जैसा दिखता है लेकिन आंखों के लिए इन छोटों के अलावा कोई अन्य चेहरे की विशेषताएं नहीं हैं। इस पोशाक में मध्यमीन शैली के काढ़े की तरह दिखने के लिए व्यवस्थित होते हैं टुकड़े से बने कवच जैसे टुकड़े शामिल हैं, यह चिपिन हड्डों से सजाया गया है, जिसमें कंधे के पैड जैसे कंधे के पैड शामिल हैं, जो उनके कंद्रे में स्टार मोटिवों से सजाए गए ढाल के रूप में ढाल के रूप में सजावट के तत्व हैं। बटन जैसे सजावटी तत्व भी हैं।



Persian: این تصویر یک اثر هنری بسیار دقیق و سورئال است که نمای نیمرخ صورت زنی را به تصویر می‌کشد و عناصر بیجیده ارگانیک در ویزگی‌های چهره او ادغام شده‌اند. پوست او با یافته‌های متنوعی شبیه برگها، تاکها و گلوهای طوفی تزئین شده که به طور یکپارچه با طرح‌های طبیعی مانند گلها و بروانه‌ها ترکیب می‌شوند. جشم او حالت اثیری دارد، با عنیبه‌های آبی درختان که با درخشش طلایی و مزه‌های بلند احاطه شده‌اند. لب‌هایش به دلیل اغراق هنری برای ایجاد اثر دراماتیک، برتر از حالت عادی هستند. ترکیب‌بندی شامل رنگ‌های زندگانی جوں طلایی، آبی، سبز، بنفش و رنگ‌های خاکی است که با لایه‌های مختلف تقاضی عمق ایجاد می‌کند.



Figure 5: **Qualitative evaluation of NEOBABEL.** Each row is based on a single concept expressed in six different languages. For clarity, we show only one of the prompts (in one language) and present six images generated from its translated prompts in the other five languages. No matter the language, NEOBABEL consistently produces semantically aligned, visually coherent results.

It supports text-guided image inpainting and extrapolation across multiple languages without requiring additional fine-tuning. As shown in Figures 6 and 7, NEOBABEL can modify or extend an input image based on prompts in different languages, producing results that remain semantically faithful and visually consistent with the adjacent visual content. These examples demonstrate the model’s ability to maintain coherence across inpainted and extrapolated regions, highlighting its potential for interactive and multilingual visual editing.

6.2.2 Cross-lingual Image Generation

A more challenging evaluation of the model’s multilingual ability involves prompts that combine multiple languages within the same input. This requires the model to integrate information from different languages into a coherent and semantically accurate image. To create these prompts, we split a base prompt into three parts and translate each into a different language. Figure 8 (in the Introduction) illustrates two such examples. The images generated by NEOBABEL demonstrate its ability to follow complex multilingual instructions, producing visually coherent and semantically faithful outputs. These results highlight the model’s cross-lingual alignment, despite not being explicitly trained for this task.

7 Ablations and Analyses

In the following sections, we conduct a series of ablation studies and analyses to evaluate the effects of progressive pretraining (Section 7.1), instruction tuning (Section 7.2), and model merging (Section 7.3). We also perform additional multilingual evaluations, including cross-linguistic consistency (Section 7.4) and code switching similarity analyses (Section 7.5).

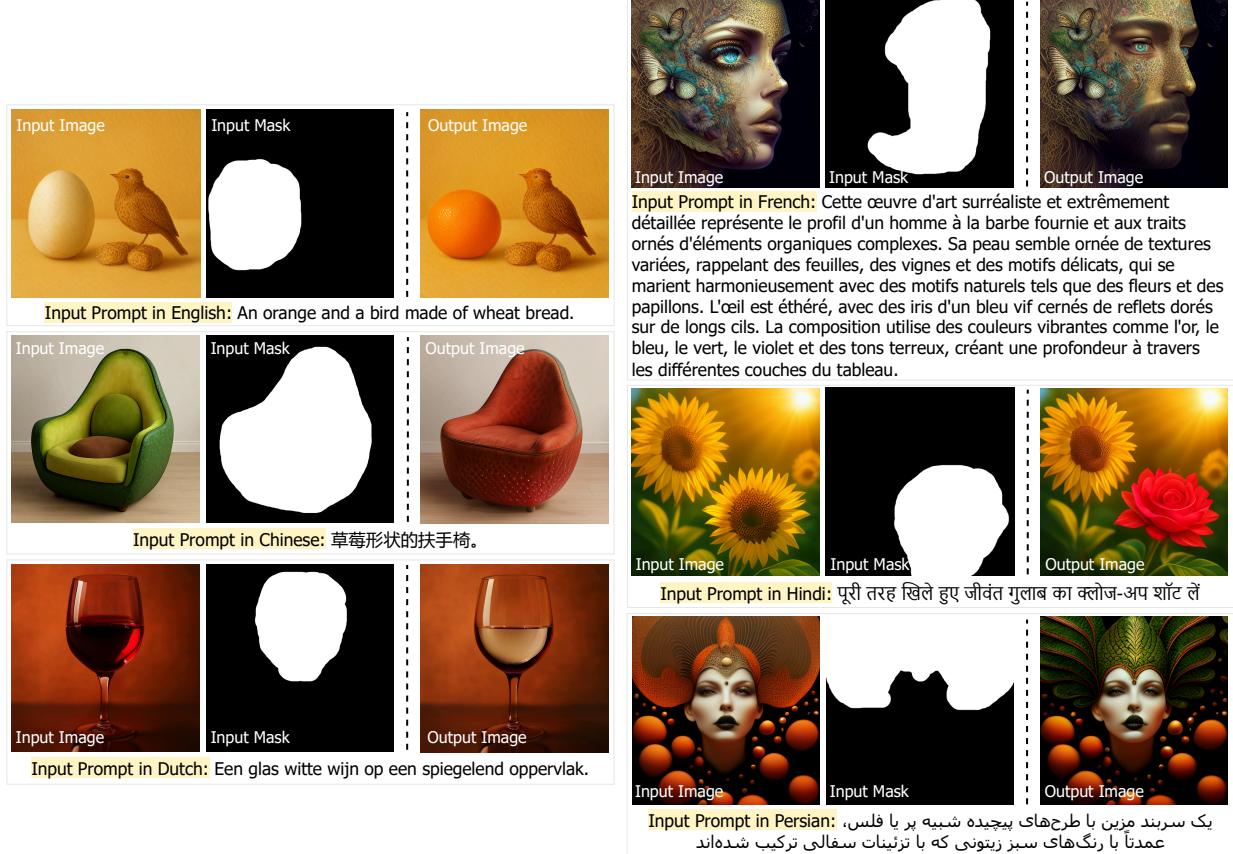


Figure 6: **Multilingual image inpainting.** NEOBABEL supports multilingual text-guided image inpainting, highlighting its potential for interactive and language-inclusive visual editing across diverse user groups.

7.1 Effect of Progressive Pretraining

We first analyze the impact of our progressive pretraining strategy across three stages at 256×256 resolution. As shown in Figure 9, each stage leads to steady improvements in multilingual performance on m-GenEval and m-DPG. In the first stage, using only m-ImageNet 1K, the average scores are modest: 0.04 on m-GenEval and 0.14 on m-DPG, indicating weak multilingual alignment. In the second stage, the addition of large-scale but noisy datasets (m-SA-1B, m-CC12M, m-LAION-Aesthetic) results in a significant increase, reaching 0.17 on m-GenEval and 0.58 on m-DPG. The average gain in performance from stage one to stage two is substantially larger on m-DPG (0.44) than on m-GenEval (0.14), suggesting that this stage improves the model’s ability to handle natural, descriptive prompts in multiple languages. The third stage incorporates higher-quality datasets (m-LAION-Aesthetic and m-JourneyDB), leading to further gains: 0.02 on m-GenEval and 0.04 on m-DPG. These results demonstrate the cumulative benefits of progressively increasing both the diversity and quality of pretraining data. While large, noisy datasets drive early generalization, high-quality data is essential for refining multilingual alignment. This staged pretraining approach provides a strong initialization for downstream instruction tuning.

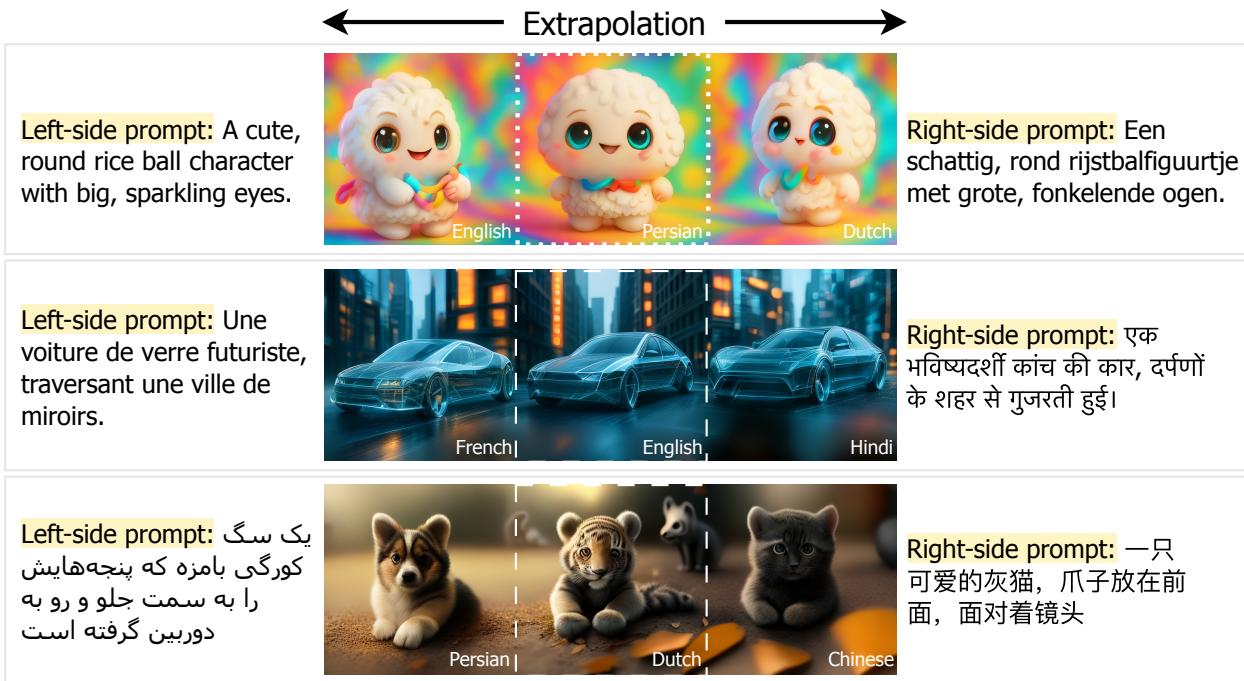


Figure 7: **Multilingual image extrapolation.** NEOBABEL successfully performs text-guided image extrapolation using multilingual prompts. Given the middle image and two different multilingual prompts (for the left and right extensions), NEOBABEL generates coherent visual completions on both sides, demonstrating extrapolation capability.

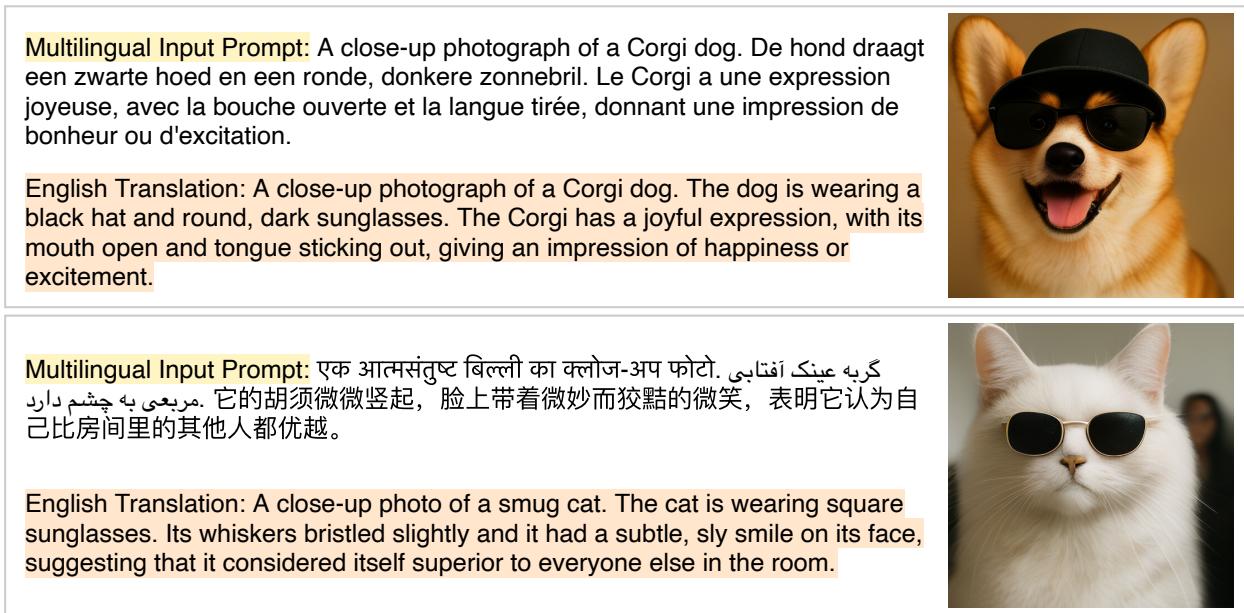


Figure 8: **Cross-Lingual Prompt Generation.** Examples of code-switched prompts mixing three languages, along with images generated by NEOBABEL. Top: English, Dutch and French. Bottom: Hindi, Persian and Chinese. English translations are shown below each prompt for reader convenience, they are not used as input.

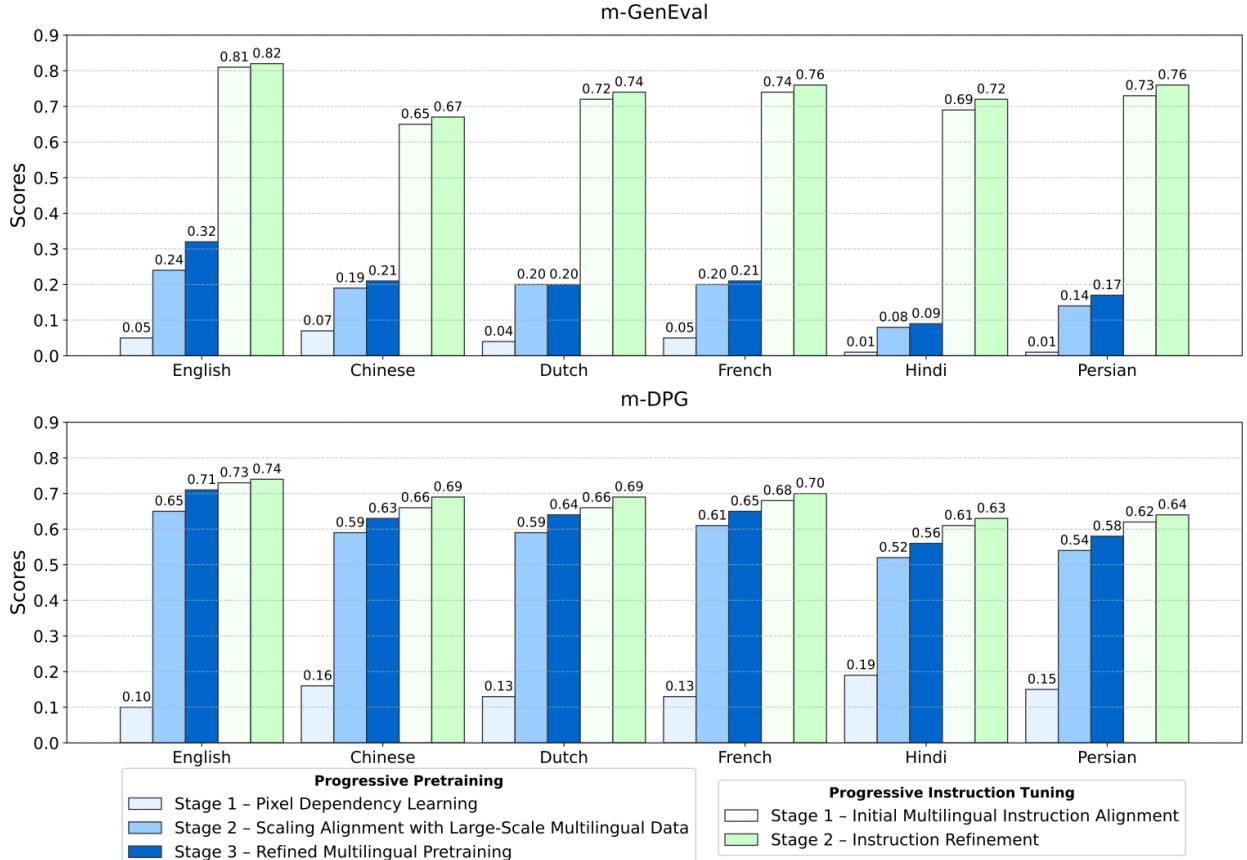


Figure 9: **Effect of Progressive Pretraining and Instruction Tuning.** Performance on m-GenEval (top) and m-DPG (bottom) improves steadily across pretraining and instruction tuning stages. Pretraining at 256×256 yields significant gains—especially on m-DPG—when scaling to large multilingual datasets (Stage 2), followed by refinement using higher-quality data (Stage 3). Instruction tuning at 512×512 brings a substantial boost (e.g., +0.52 on m-GenEval) and further improvement by increasing the share of curated, instruction-aligned samples. These results highlight how scale drives broad generalization, while data quality and resolution are key for performance refinement.

7.2 Effect of Progressive Instruction Tuning

We analyze the effect of high-resolution instruction tuning using a fixed dataset mixture: m-LAION-Aesthetic, m-JourneyDB, and m-BLIP3o-Instruct, all at 512×512 resolution. As shown in Figure 9, both tuning stages progressively refine multilingual alignment, with consistent gains across all six languages. In the first stage, each training batch consists of 60% m-LAION-Aesthetic, 30% m-JourneyDB, and 10% m-BLIP3o-Instruct samples. This stage yields a substantial multilingual gain of 0.52 on m-GenEval and 0.04 on m-DPG compared to the final stage of pretraining, indicating that high-resolution supervision provides strong improvements when combined with highly curated instruction tuning dataset. In the second stage, each training batch shifts emphasis toward higher-quality and instruction-aligned data, consisting of 25% m-LAION-Aesthetic, 60% m-JourneyDB, and 15% m-BLIP3o-Instruct samples. This leads to a further multilingual gain of 0.02 in m-GenEval and a boost in m-DPG. These results show that beyond increasing the resolution, the relative weight of curated and instruction-focused datasets plays a pivotal role in shaping multilingual capability. Prioritizing high-quality supervision at higher resolution proves effective for achieving competitive alignment ultimately enabling our 2B model to rival much larger models.

7.3 Effect of Model Merging on Generalization

We investigate the impact of model merging on multilingual image generation performance by combining $N = 20$ checkpoints sampled at 10,000-step intervals from the second instruction tuning stage (steps 0–200K). Table 4 reports the results of three merging strategies: Simple Moving Average (SMA), Exponential Moving Average (EMA), and Weighted Moving Average (WMA) compared to the last checkpoint baseline. As reported, m-GenEval score on English prompt improves from 0.81 to 0.83 after model merging. Both WMA and SMA reach this upper bound, indicating that merging checkpoints along the optimization path enhances semantic alignment. Moreover, m-DPG score on English prompt remains stable or show modest gains, suggesting that merging preserves the model’s ability to accurately follow dense, attribute-rich prompts without sacrificing fine-grained multilingual grounding. Among the merging strategies, SMA performs best overall due to its uniform averaging over well-aligned checkpoints. EMA also improves results but remains more susceptible to short-term training noise. WMA offers a compromise by emphasizing later checkpoints, trading off stability for adaptability. These findings underscore that checkpoint merging can meaningfully enhance both compositional understanding and multilingual robustness, with SMA offering a simple yet effective strategy.

7.4 Cross-Lingual Consistency Analysis

This ablation examines how different models preserve consistency in image generation across languages, evaluated using the CLC score introduced in Section 5.2 with both EVA-CLIP and DINOv2 vision encoders (Table 5).

Across both backbones, NEOBABEL achieves the highest scores, 0.79 (EVA-CLIP) and 0.61 (DINOv2), outperforming larger models such as BLIP3-o 8B (0.77/0.45) and Janus Pro 7B (0.67/0.30).

This indicates that training strategy and data alignment play a more significant role than parameter count alone in achieving cross-lingual consistency. In the EVA-CLIP based CLC, which emphasizes high-level semantic similarity due to its contrastive training objective with text, high scores reflect strong cross-lingual consistency in scene-level concept. In contrast, DINOv2-based CLC captures visual-structural coherence making it more sensitive to differ-

Method	m-GenEval	m-DPG
Last checkpoint	0.81	0.73
EMA	0.82	0.75
WMA	0.83	0.75
SMA	0.83	0.75

Table 4: **Effect of model merging on generalization.** Without any fine-tuning, all merging strategies improve performance on m-GenEval and slightly enhance m-DPG, highlighting model merging as a simple yet effective way to boost generalization. This ablation uses English prompts.

Both WMA and SMA reach this upper bound, indicating that merging checkpoints along the optimization path enhances semantic alignment. Moreover, m-DPG score on English prompt remains stable or show modest gains, suggesting that merging preserves the model’s ability to accurately follow dense, attribute-rich prompts without sacrificing fine-grained multilingual grounding. Among the merging strategies, SMA performs best overall due to its uniform averaging over well-aligned checkpoints. EMA also improves results but remains more susceptible to short-term training noise. WMA offers a compromise by emphasizing later checkpoints, trading off stability for adaptability. These findings underscore that checkpoint merging can meaningfully enhance both compositional understanding and multilingual robustness, with SMA offering a simple yet effective strategy.

Model	Params	EVA-CLIP	DINOv2
Show-o	1.3B	0.47	0.16
Janus	1.3B	0.67	0.26
Janus Pro	7B	0.67	0.30
BLIP3-o	4B	0.76	0.44
BLIP3-o	8B	0.77	0.45
NEOBABEL	2B	0.79	0.61

Table 5: **Cross-Lingual Consistency Analysis** using CLC scores with EVA-CLIP and DINOv2 backbones. NEOBABEL (2B) outperforms larger models, showing stronger cross-lingual consistency in both semantic and visual domains. The larger gap in DINOv2 scores highlights its greater sensitivity to visual-structural variations, revealing inconsistencies that EVA-CLIP’s semantic focus may overlook.

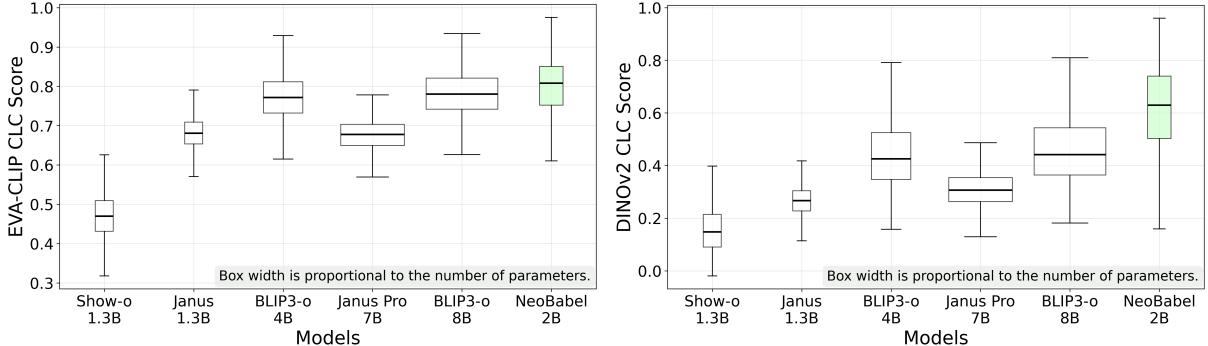


Figure 10: **Cross-Lingual Consistency (CLC) Score Distributions across Models.** We show the distribution of CLC scores computed using EVA-CLIP (left column) and DINOV2 (right column), where higher values reflect greater consistency across languages. EVA-CLIP captures semantic similarity, while DINOV2 is more sensitive to visual structure and layout. NEOBABEL achieves the highest scores under both backbones—particularly with DINOV2—demonstrating strong alignment in both meaning and visual form. Box widths reflect model size, showing that bigger models aren’t always more consistent.

ences in object composition, layout, or fine-grained visual patterns. Notably, the relative margin between NEOBABEL and other models is more pronounced under DINOV2. For instance, BLIP3-o’s DINOV2 score drops significantly despite competitive EVA-CLIP CLC score suggesting its outputs vary more in visual structure across languages.

This amplifies the role of multilingual alignment not just in semantics, but in visual form a dimension better captured by DINOV2. NEOBABEL achieves high scores under both backbones, thereby demonstrating strong cross-lingual consistency in both scene semantics and visual structure. Figure 10 complements Table 5 by visualizing the distribution of CLC scores for each model using both EVA-CLIP and DINOV2 backbones. In terms of EVA-CLIP based CLC variation, NEOBABEL performs on par with BLIP3-o 8B, the second-best model. However, under the DINOV2-based CLC variation, NEOBABEL outperforms all baselines, exhibiting lower dispersion and higher consistency. We can also observe from the figure that simply increasing model size does not guarantee better cross-lingual consistency, as evidenced by the lower DINOV2 scores of Janus Pro 7B and BLIP3-o 8B compared to NEOBABEL.

7.5 Code Switching Similarity Analysis

We evaluate model robustness to intra-prompt code switching using the proposed CSS score in Section 5.2 with EVA-CLIP and DINOV2 backbones. As shown in Table 6 and Figure 11, NEOBA-

Model	Params	EVA-CLIP		DINOv2	
		EF	ES	EF	ES
Show-o	1.3B	0.73	0.72	0.41	0.38
Janus	1.3B	0.75	0.73	0.50	0.43
Janus Pro	7B	0.76	0.72	0.58	0.50
BLIP3-o	4B	0.75	0.75	0.54	0.54
BLIP3-o	8B	0.74	0.74	0.52	0.51
NEOBABEL	2B	0.82	0.81	0.67	0.64

Table 6: **Code Switching Similarity (CSS) Analysis** using EVA-CLIP and DINOV2 backbones. Scores are reported for two prompt variants: English First (EF) and English Second (ES). NEOBABEL (2B) outperforms larger models, showing strong visual consistency and robustness to code-mixed input order. The larger DINOV2 gap reflects its higher sensitivity to visual-structural variation, while EVA-CLIP remains more stable due to its semantic focus.

BEL consistently outperforms larger models, demonstrating stronger visual consistency under both English-First (EF) and English-Second (ES) variations. Under EVA-CLIP, all models exhibit minimal difference between EF and ES, suggesting that the position of the English segment has limited effect on global semantic alignment. In contrast, DINOv2 scores are lower across the board, indicating greater difficulty in maintaining consistent visual structure when mixing languages.

Importantly, a desirable outcome is both high CSS scores (indicating alignment with the reference image) and minimal gap between EF and ES (indicating robustness to code-switch position). NEOBABEL achieves this balance, with CSS scores of 0.82 (EF) and 0.81 (ES) for EVA-CLIP, and 0.67 (EF) and 0.64 (ES) for DINOv2. The box plots reveal that although larger models like BLIP3-o (8B) achieve competitive means, they show greater variability across prompts. NEOBABEL demonstrates both higher median performance and lower dispersion, confirming its consistent handling of code-mixed inputs. These results further highlight that scaling model size does not necessarily improve code-mixed prompt robustness effective multilingual alignment plays a larger role.

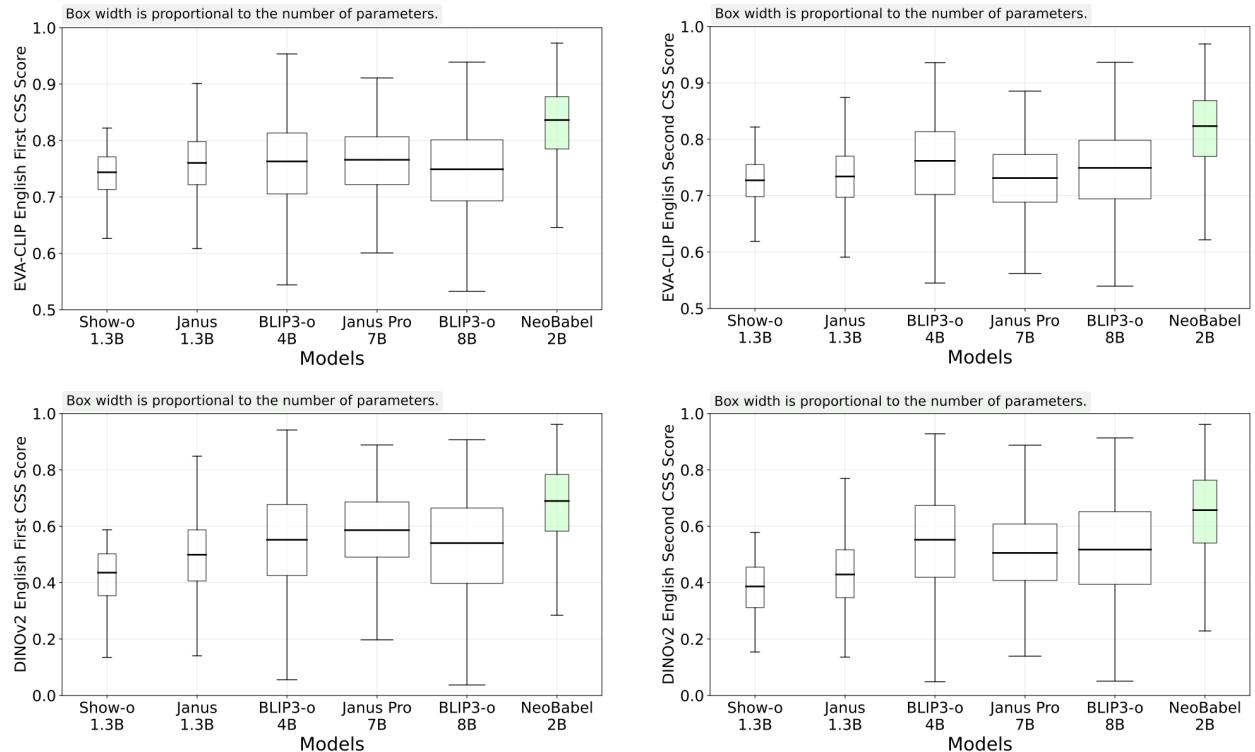


Figure 11: **Variation in Code Switching Similarity (CSS) Scores across Models.** We report CSS scores for code-mixed prompts under two settings: English-first (left column) and English-second (right column), using EVA-CLIP (top row) and DINOv2 (bottom row) as backbones. Higher scores indicate stronger visual alignment with the reference image, while smaller EF–ES gaps suggest robustness to code-switch position. NEOBABEL consistently achieves higher medians and lower variance than larger baselines, especially under DINOv2, highlighting its effective and stable handling of multilingual prompts.

8 Related Works

Large Multimodal Models. Recent advances in large multimodal models (LMMs) [Liu et al., 2024b; Chen et al., 2024b; Li et al., 2024a; Bai et al., 2025; Dash et al., 2025] have extended large language models (LLMs) [Touvron et al., 2023; Yang et al., 2024] to support image understanding

tasks, including image captioning and visual question answering. These models typically rely on a vision encoder to extract image features, which are then projected into the LLM embedding space for cross-modal alignment. More recent encoder-free models [Xie et al., 2025b; Diao et al., 2024; 2025] bypass the explicit image encoder and instead align raw visual tokens directly within the LLM space. Among early efforts to enable multilingual visual understanding, Maya [Alam et al., 2024; 2025], Aya-Vision [Dash et al., 2025] and Pangea [Yue et al., 2024] incorporate a multilingual training corpus. However, they are limited to image understanding tasks. In contrast, our proposed NEOBABEL architecture focuses exclusively on multilingual image generation, offering the first encoder-free model that aligns visual features in the LLM space while supporting cross-lingual generation. Architecturally, NEOBABEL is closely related to show-o [Xie et al., 2025b], sharing the same design goal of direct visual alignment in language space, but differs in its task focus and multilingual design.

Visual Generative Models. Two dominant paradigms have emerged for image and video generation: diffusion-based [Rombach et al., 2022; Peebles & Xie, 2023; Bao et al., 2023; Chen et al., 2024a; Xie et al., 2023; Wu et al., 2023a; Lipman et al., 2022; Xie et al., 2025a; Qin et al., 2025; Zhang et al., 2023; Seawead et al., 2025] and autoregressive [Sun et al., 2024; Kondratyuk et al., 2023; Chen et al., 2020; Pang et al., 2024; Li et al., 2025a] models. Diffusion models typically combine pretrained text encoders with denoising networks to iteratively refine visual outputs, while autoregressive models adopt LLM-based architectures trained via next-token prediction. Recent hybrid approaches [Li et al., 2024b; Liu et al., 2024c; Fan et al., 2025] attempt to unify the strengths of both paradigms for more powerful generation. NEOBABEL follows the diffusion-based paradigm but distinguishes itself by adopting an LLM-style architecture for visual token modeling. This removes the reliance on frozen text encoders and instead builds on top of a strong multilingual decoder-based LLM, enabling tighter integration between language and vision.

Unified Multimodal Models. Unified multimodal models aim to handle both image understanding and generation within a single architecture, typically categorized into native and adapter-based approaches. Native approaches such as Chameleon [Team, 2024], Show-o [Xie et al., 2025b], and Transfusion [Zhou et al., 2025] adopt either autoregressive, diffusion, or hybrid modeling strategies to jointly process vision and language. Recent work [Wang et al., 2024; Wu et al., 2024; Ma et al., 2025; Jiao et al., 2025; Chen et al., 2025c; Song et al., 2025] has focused on improving tokenization and training efficiency to enhance cross-modal alignment. A parallel direction [Tang et al., 2023; Lu et al., 2023; Dong et al., 2024; Ge et al., 2024; Tong et al., 2024; Pan et al., 2025; Chen et al., 2025a; Wu et al., 2023b] constructs unified multimodal models by connecting pretrained LMMs and generative models via adapters or learnable tokens. While modular and flexible, these systems often rely on frozen components and lack full cross-modal integration. Our model, NEOBABEL, aligns more closely with native unified multimodal models by unifying visual and textual modeling within a single decoder-based architecture, without relying on adapters or frozen backbones. Although NEOBABEL supports multilingual multimodal understanding, this work focuses specifically on multilingual image generation.

9 Limitations

While NEOBABEL demonstrates strong multilingual image generation capabilities, several limitations remain. First, the model currently supports only six languages; extending to broader linguistic coverage would require further tokenizer adaptation and additional training. Second, although

NEOBABEL adopts a unified architecture, it does not yet support vision-language tasks such as visual question answering, due to the absence of task-specific fine-tuning. Third, the model’s performance is constrained by its parameter size and the diversity and quality of the training data. For instance, during instruction tuning, we used a fixed mixture of m-LAION-Aesthetic, m-JourneyDB, and m-BLIP3o-Instruct, without performing an extensive sweep over mixture ratios—an area that could reveal further improvements. We leave these directions, including task expansion, larger-scale scaling, and wider language coverage, for future research.

10 Conclusion

NEOBABEL demonstrates that high-quality, efficient multilingual image generation is not only possible but also advantageous. Through strategic data curation and a unified architecture, we set a new Pareto frontier in performance, efficiency, and inclusivity. While currently focused on text-to-image generation, our model is structurally capable of broader multimodal tasks. Our results across m-GenEval and m-DPG benchmarks, paired with the introduction of new evaluation metrics (CLC and CSS), establish a robust foundation for the next generation of multilingual generative models.

Our work opens several promising avenues for future research. First, extending NEOBABEL to encompass a wider variety of languages, particularly those currently underrepresented in vision-language research, remains an important objective. The modularity of our framework and the accompanying open-source toolkit are designed to facilitate such extensions. Second, beyond linguistic diversity, integrating cultural grounding into multimodal models presents a compelling research direction. By curating datasets annotated with region-specific concepts, aesthetic preferences, and social norms, future work could develop models that are not only multilingual but also culturally aware and adaptive. Finally, this work contributes to the broader goal of democratizing generative AI. By releasing all model weights, datasets, and evaluation protocols, we aim to encourage the research community to build upon this foundation, ultimately advancing toward generative models that better reflect and serve global linguistic and cultural diversity.

Acknowledgment

We would like to thank the Cohere Labs team for their valuable feedback and for providing generous computing resources for conducting and analyzing our experiments. We further acknowledge the Dutch Research Council (NWO) in The Netherlands for awarding this project access to the LUMI supercomputer, owned by the EuroHPC Joint Undertaking, hosted by CSC (Finland) and the LUMI consortium through the Computing Time on National Computer Facilities call. We also acknowledge NWO for providing access to Snellius, hosted by SURF through the Computing Time on National Computer Facilities call for proposals. Cees G. M. Snoek is (partially) funded by the Horizon Europe project ELLIOT (GA No. 101214398).

References

Nahid Alam, Karthik Reddy Kanjula, Surya Guthikonda, Timothy Chung, Bala Krishna S Vegesna, Abhipsha Das, Anthony Susevski, Ryan Sze-Yin Chan, SM Uddin, Shayekh Bin Islam, et al. Maya:

An instruction finetuned multilingual multimodal model. *arXiv preprint arXiv:2412.07112*, 2024.

Nahid Alam, Karthik Reddy Kanjula, Surya Guthikonda, Timothy Chung, Bala Krishna S Vigesna, Abhipsha Das, Anthony Susevski, Ryan Sze-Yin Chan, SM Uddin, Shayekh Bin Islam, et al. Behind maya: Building a multilingual vision language model. *arXiv preprint arXiv:2505.08910*, 2025.

Niyati Bafna, Tianjian Li, Kenton Murray, David R Mortensen, David Yarowsky, Hale Sirin, and Daniel Khashabi. The translation barrier hypothesis: Multilingual generation with large language models suffers from implicit translation failure. *arXiv preprint arXiv:2506.22724*, 2025.

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.

Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth words: A vit backbone for diffusion models. In *Conference on Computer Vision and Pattern Recognition*, 2023.

Elisa Bassignana, Amanda Cercas Curry, and Dirk Hovy. The ai gap: How socioeconomic status affects language technology interactions. *arXiv preprint arXiv:2505.12158*, 2025.

Lisa Beinborn and Rochelle Choenni. Semantic drift in multilingual representations. *Computational Linguistics*, 2020.

Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *Conference on Computer Vision and Pattern Recognition*, 2021.

Juhai Chen, Zhiyang Xu, Xichen Pan, Yushi Hu, Can Qin, Tom Goldstein, Lifu Huang, Tianyi Zhou, Saining Xie, Silvio Savarese, et al. Blip3-o: A family of fully open unified multimodal models-architecture, training and dataset. *arXiv preprint arXiv:2505.09568*, 2025a.

Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Zhongdao Wang, James T. Kwok, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart- α : Fast training of diffusion transformer for photorealistic text-to-image synthesis. In *International Conference on Learning Representations*, 2024a.

Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pretraining from pixels. In *International Conference on Machine Learning*, 2020.

Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan. Janus-pro: Unified multimodal understanding and generation with data and model scaling. *arXiv preprint arXiv:2501.17811*, 2025b.

Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024b.

Zhe Chen, Jiannan Wu, Wenhui Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Conference on Computer Vision and Pattern Recognition*, 2024c.

Zisheng Chen, Chunwei Wang, Xiuwei Chen, Hang Xu, Jianhua Han, and Xiaodan Liang. Semhitok: A unified image tokenizer via semantic-guided hierarchical codebook for multimodal understanding and generation. *arXiv preprint arXiv:2503.06764*, 2025c.

Reuben Cohn-Gordon and Noah Goodman. Lost in machine translation: A method to reduce meaning loss. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2019.

Marta R Costa-Jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*, 2022.

Saurabh Dash, Yiyang Nan, John Dang, Arash Ahmadian, Shivalika Singh, Madeline Smith, Bharat Venkitesh, Vlad Shmyhlo, Viraat Aryabumi, Walter Beller-Morales, Jeremy Pekmez, Jason Ozuzu, Pierre Richemond, Acyr Locatelli, Nick Frosst, Phil Blunsom, Aidan Gomez, Ivan Zhang, Marzieh Fadaee, Manoj Govindassamy, Sudip Roy, Matthias Gallé, Beyza Ermis, Ahmet Üstün, and Sara Hooker. Aya vision: Advancing the frontier of multilingual multimodality. *arXiv preprint arXiv:2505.08751*, 2025.

Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadmreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, et al. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. In *Conference on Computer Vision and Pattern Recognition*, 2025.

Haiwen Diao, Yufeng Cui, Xiaotong Li, Yueze Wang, Huchuan Lu, and Xinlong Wang. Unveiling encoder-free vision-language models. *arXiv preprint arXiv:2406.11832*, 2024.

Haiwen Diao, Xiaotong Li, Yufeng Cui, Yueze Wang, Haoge Deng, Ting Pan, Wenxuan Wang, Huchuan Lu, and Xinlong Wang. Evev2: Improved baselines for encoder-free vision-language models. *arXiv preprint arXiv:2502.06788*, 2025.

Runpei Dong, Chunrui Han, Yuang Peng, Zekun Qi, Zheng Ge, Jinrong Yang, Liang Zhao, Jianjian Sun, Hongyu Zhou, Haoran Wei, Xiangwen Kong, Xiangyu Zhang, Kaisheng Ma, and Li Yi. DreamLLM: Synergistic multimodal comprehension and creation. In *International Conference on Learning Representations*, 2024.

Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *International Conference on Machine Learning*, 2024.

Fahim Faisal and Antonios Anastasopoulos. An efficient approach for studying cross-lingual transfer in multilingual language models. In Jonne Sälevä and Abraham Owodunni (eds.), *Proceedings of the Fourth Workshop on Multilingual Representation Learning*, 2024.

Lijie Fan, Luming Tang, Siyang Qin, Tianhong Li, Xuan Yang, Siyuan Qiao, Andreas Steiner, Chen Sun, Yuanzhen Li, Tao Zhu, et al. Unified autoregressive visual generation and understanding with continuous tokens. *arXiv preprint arXiv:2503.13436*, 2025.

Felix Friedrich, Katharina Hammerl, Patrick Schramowski, Manuel Brack, Jindrich Libovicky, Christian Kersting, and Alexander Fraser. Multilingual text-to-image generation magnifies gender stereotypes and prompt engineering may not help you. *arXiv preprint arXiv:2401.16092*, 2024.

Yuying Ge, Sijie Zhao, Jinguo Zhu, Yixiao Ge, Kun Yi, Lin Song, Chen Li, Xiaohan Ding, and Ying Shan. Seed-x: Multimodal models with unified multi-granularity comprehension and generation. *arXiv preprint arXiv:2404.14396*, 2024.

Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.

Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework for evaluating text-to-image alignment. *Advances on Neural Information Processing Systems*, 2023.

Alex Henry, Prudhvi Raj Dachapally, Shubham Pawar, and Yuxuan Chen. Query-key normalization for transformers. *arXiv preprint arXiv:2010.04245*, 2020.

Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.

Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip diffusion models with llm for enhanced semantic alignment. *arXiv preprint arXiv:2403.05135*, 2024.

Shaoxiong Ji, Timothee Mickus, Vincent Segonne, and Jörg Tiedemann. Can machine translation bridge multilingual pretraining and cross-lingual transfer learning? *arXiv preprint arXiv:2403.16777*, 2024.

Yang Jiao, Haibo Qiu, Zequn Jie, Shaoxiang Chen, Jingjing Chen, Lin Ma, and Yu-Gang Jiang. Unitoken: Harmonizing multimodal understanding and generation through unified visual encoding. *arXiv preprint arXiv:2504.04423*, 2025.

Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, Hervé Jégou, and Tomas Mikolov. Fasttext.zip: Compressing text classification models. *arXiv preprint arXiv:1612.03651*, 2016.

Nithish Kannen, Arif Ahmad, marco Andreetto, Vinodkumar Prabhakaran, Utsav Prabhu, Adji Bouocco Dieng, Pushpak Bhattacharyya, and Shachi Dave. Beyond aesthetics: Cultural competence in text-to-image models, 2024.

Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *IEEE International Conference on Computer Vision*, 2023.

Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, Jonathan Huang, Rachel Hornung, Hartwig Adam, Hassan Akbari, Yair Alon, Vighnesh Birodkar, et al. Videopoet: A large language model for zero-shot video generation. *arXiv preprint arXiv:2312.14125*, 2023.

Julia Kreutzer, Eleftheria Briakou, Sweta Agrawal, Marzieh Fadaee, and Kocmi Tom. Déjà vu: Multilingual llm evaluation through the lens of machine translation evaluation, 2025.

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.

Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024a.

Haopeng Li, Jinyue Yang, Guoqi Li, and Huan Wang. Autoregressive image generation with randomized parallel decoding. *arXiv preprint arXiv:2503.10568*, 2025a.

Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image generation without vector quantization. *arXiv preprint arXiv:2406.11838*, 2024b.

Yafu Li, Ronghao Zhang, Zhilin Wang, Huajian Zhang, Leyang Cui, Yongjing Yin, Tong Xiao, and Yue Zhang. Lost in literalism: How supervised training shapes translationese in llms, 2025b.

Yunshui Li, Yiyuan Ma, Shen Yan, Chaoyi Zhang, Jing Liu, Jianqiao Lu, Ziwen Xu, Mengzhao Chen, Minrui Wang, Shiyi Zhan, et al. Model merging in pre-training of large language models. *arXiv preprint arXiv:2505.12082*, 2025c.

Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.

Bingshuai Liu, Longyue Wang, Chenyang Lyu, Yong Zhang, Jinsong Su, Shuming Shi, and Zhaopeng Tu. On the cultural gap in text-to-image generation, 2023.

Hao Liu, Wilson Yan, Matei Zaharia, and Pieter Abbeel. World model on million-length video and language with blockwise ringattention. *arXiv preprint arXiv:2402.08268*, 2024a.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances on Neural Information Processing Systems*, 2024b.

Haozhe Liu, Shikun Liu, Zijian Zhou, Mengmeng Xu, Yanping Xie, Xiao Han, Juan C. Pérez, Ding Liu, Kumara Kahatapitiya, Menglin Jia, Jui-Chieh Wu, Sen He, Tao Xiang, Jürgen Schmidhuber, and Juan-Manuel Pérez-Rúa. Mardini: Masked autoregressive diffusion for video generation at scale. *arXiv preprint arXiv:2410.20280*, 2024c.

Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha Kembhavi. Unified-io 2: Scaling autoregressive multimodal models with vision, language, audio, and action. *arXiv preprint arXiv:2312.17172*, 2023.

Chuofan Ma, Yi Jiang, Junfeng Wu, Jihan Yang, Xin Yu, Zehuan Yuan, Bingyue Peng, and Xiaojuan Qi. Unitok: A unified tokenizer for visual generation and understanding. *arXiv preprint arXiv:2502.20321*, 2025.

Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.

Xichen Pan, Satya Narayan Shukla, Aashu Singh, Zhuokai Zhao, Shlok Kumar Mishra, Jialiang Wang, Zhiyang Xu, Jiahui Chen, Kunpeng Li, Felix Juefei-Xu, Ji Hou, and Saining Xie. Transfer between modalities with metaqueries. *arXiv preprint arXiv:2504.06256*, 2025.

Ziqi Pang, Tianyuan Zhang, Fujun Luan, Yunze Man, Hao Tan, Kai Zhang, William T. Freeman, and Yu-Xiong Wang. Randar: Decoder-only autoregressive visual generation in random orders. *arXiv preprint arXiv:2412.01827*, 2024.

William Peebles and Saining Xie. Scalable diffusion models with transformers. In *IEEE International Conference on Computer Vision*, 2023.

Aidan Peppin, Julia Kreutzer, Alice Schoenauer Sebag, Kelly Marchisio, Beyza Ermis, John Dang, Samuel Cahyawijaya, Shivalika Singh, Seraphina Goldfarb-Tarrant, Viraat Aryabumi, Aakanksha, Wei-Yin Ko, Ahmet Üstün, Matthias Gallé, Marzieh Fadaee, and Sara Hooker. The multilingual divide and its impact on global ai safety, 2025.

Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023.

Qi Qin, Le Zhuo, Yi Xin, Ruoyi Du, Zhen Li, Bin Fu, Yiting Lu, Jiakang Yuan, Xinyue Li, Dongyang Liu, et al. Lumina-image 2.0: A unified and efficient image generative framework. *arXiv preprint arXiv:2503.21758*, 2025.

Liao Qu, Huichao Zhang, Yiheng Liu, Xu Wang, Yi Jiang, Yiming Gao, Hu Ye, Daniel K Du, Zehuan Yuan, and Xinglong Wu. Tokenflow: Unified image tokenizer for multimodal understanding and generation. In *Conference on Computer Vision and Pattern Recognition*, 2025.

Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.

Aniket Rege, Zinnia Nie, Mahesh Ramesh, Unmesh Raskar, Zhuoran Yu, Aditya Kusupati, Yong Jae Lee, and Ramya Korlakai Vinayak. Cure: Cultural gaps in the long tail of text-to-image systems, 2025.

Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Conference on Computer Vision and Pattern Recognition*, 2022.

Israfel Salazar, Manuel Fernández Burda, Shayekh Bin Islam, Arshia Soltani Moakhar, Shivalika Singh, Fabian Farestam, Angelika Romanou, Danylo Boiko, Dipika Khullar, Mike Zhang, et al. Kaleidoscope: In-language exams for massively multilingual vision evaluation. *arXiv preprint arXiv:2504.07072*, 2025.

Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances on Neural Information Processing Systems*, 2022.

Team Seawead, Ceyuan Yang, Zhijie Lin, Yang Zhao, Shanchuan Lin, Zhibei Ma, Haoyuan Guo, Hao Chen, Lu Qi, Sen Wang, et al. Seaweed-7b: Cost-effective training of video generation foundation model. *arXiv preprint arXiv:2504.08685*, 2025.

Luisa Shimabucoro, Ahmet Ustun, Marzieh Fadaee, and Sebastian Ruder. A post-trainer’s guide to multilingual training data: Uncovering cross-lingual transfer dynamics. *arXiv preprint arXiv:2504.16677*, 2025.

Shivalika Singh, Angelika Romanou, Clémentine Fourrier, David I Adelani, Jian Gang Ngui, Daniel Vila-Suero, Peerat Limkonchotiwat, Kelly Marchisio, Wei Qi Leong, Yosephine Susanto, et al. Global mmlu: Understanding and addressing cultural and linguistic biases in multilingual evaluation. *arXiv preprint arXiv:2412.03304*, 2024.

Wei Song, Yuran Wang, Zijia Song, Yadong Li, Haoze Sun, Weipeng Chen, Zenan Zhou, Jianhua Xu, Jiaqi Wang, and Kaicheng Yu. Dualtoken: Towards unifying visual understanding and generation with dual visual vocabularies. *arXiv preprint arXiv:2503.14324*, 2025.

Keqiang Sun, Junting Pan, Yuying Ge, Hao Li, Haodong Duan, Xiaoshi Wu, Renrui Zhang, Aojun Zhou, Zipeng Qin, Yi Wang, et al. Journeydb: A benchmark for generative image understanding. *Advances on Neural Information Processing Systems*, 2023a.

Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan. Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint arXiv:2406.06525*, 2024.

Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training techniques for clip at scale. *arXiv preprint arXiv:2303.15389*, 2023b.

Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. Any-to-any generation via composable diffusion. *Advances on Neural Information Processing Systems*, 2023.

Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024.

Shengbang Tong, David Fan, Jiachen Zhu, Yunyang Xiong, Xinlei Chen, Koustuv Sinha, Michael Rabbat, Yann LeCun, Saining Xie, and Zhuang Liu. Metamorph: Multimodal understanding and generation via instruction tuning. *arXiv preprint arXiv:2412.14164*, 2024.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *Computing Research Repository*, 2023.

Eva Vanmassenhove, Dimitar Shterionov, and Andy Way. Lost in translation: Loss and decay of linguistic richness in machine translation. In *Proceedings of Machine Translation Summit XVII: Research Track*, pp. 222–232, Dublin, Ireland, 2019.

Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, Quan Sun, Yufeng Cui, Jinsheng Wang, Fan Zhang, Yueze Wang, Zhen Li, Qiying Yu, et al. Emu3: Next-token prediction is all you need. *arXiv preprint arXiv:2409.18869*, 2024.

Shira Wein and Nathan Schneider. Lost in translationese? reducing translation effect using abstract meaning representation. *arXiv preprint arXiv:2304.11501*, 2023.

Chengyue Wu, Xiaokang Chen, Zhiyu Wu, Yiyang Ma, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, Chong Ruan, et al. Janus: Decoupling visual encoding for unified multimodal understanding and generation. In *Conference on Computer Vision and Pattern Recognition*, 2025.

Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion

models for text-to-video generation. In *IEEE International Conference on Computer Vision*, 2023a.

Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. Next-gpt: Any-to-any multi-modal llm. *arXiv preprint arXiv:2309.05519*, 2023b.

Yecheng Wu, Zhuoyang Zhang, Junyu Chen, Haotian Tang, Dacheng Li, Yunhao Fang, Ligeng Zhu, Enze Xie, Hongxu Yin, Li Yi, et al. Vila-u: a unified foundation model integrating visual understanding and generation. *arXiv preprint arXiv:2409.04429*, 2024.

Enze Xie, Junsong Chen, Yuyang Zhao, Jincheng Yu, Ligeng Zhu, Chengyue Wu, Yujun Lin, Zhekai Zhang, Muyang Li, Junyu Chen, et al. Sana 1.5: Efficient scaling of training-time and inference-time compute in linear diffusion transformer. *arXiv preprint arXiv:2501.18427*, 2025a.

Jinheng Xie, Yuexiang Li, Yawen Huang, Haozhe Liu, Wentian Zhang, Yefeng Zheng, and Mike Zheng Shou. Boxdiff: Text-to-image synthesis with training-free box-constrained diffusion. In *IEEE International Conference on Computer Vision*, 2023.

Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer to unify multimodal understanding and generation. *International Conference on Learning Representations*, 2025b.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.

Lijun Yu, José Lezama, Nitesh B Gundavarapu, Luca Versari, Kihyuk Sohn, David Minnen, Yong Cheng, Vighnesh Birodkar, Agrim Gupta, Xiuye Gu, et al. Language model beats diffusion–tokenizer is key to visual generation. *arXiv preprint arXiv:2310.05737*, 2023.

Xiang Yue, Yueqi Song, Akari Asai, Seungone Kim, Jean de Dieu Nyandwi, Simran Khanuja, Anjali Kantharuban, Lintang Sutawika, Sathyanarayanan Ramamoorthy, and Graham Neubig. Pangea: A fully open multilingual multimodal llm for 39 languages. In *International Conference on Learning Representations*, 2024.

David Junhao Zhang, Jay Zhangjie Wu, Jia-Wei Liu, Rui Zhao, Lingmin Ran, Yuchao Gu, Difei Gao, and Mike Zheng Shou. Show-1: Marrying pixel and latent diffusion models for text-to-video generation. *arXiv preprint arXiv:2309.15818*, 2023.

Chunting Zhou, LILI YU, Arun Babu, Kushal Tirumala, Michihiro Yasunaga, Leonid Shamis, Jacob Kahn, Xuezhe Ma, Luke Zettlemoyer, and Omer Levy. Transfusion: Predict the next token and diffuse images with one multi-modal model. In *International Conference on Learning Representations*, 2025.

Appendix A

This appendix provides additional training details and qualitative results to supplement the main paper. Table 7 outlines the key hyperparameters used across the three pretraining stages and two instruction tuning stages of NEOBABEL. Figure 12 presents representative multilingual generation examples, demonstrating visual consistency across six languages.

Hyperparameters	Pretraining			Instruction Tuning	
	1st Stage	2nd Stage	3rd Stage	1st Stage	2nd Stage
Training Steps	500k	500k	500k	500k	200k
Warmup Steps	5000	5000	5000	5000	2000
Learning Rate	$1e - 4$	$1e - 4$	$1e - 4$	$2e - 4$	$5e - 05$
Learning Rate Decay	cosine	cosine	cosine	cosine	cosine
Optimizer	AdamW	AdamW	AdamW	AdamW	AdamW
Image Resolution	256×256	256×256	256×256	512×512	512×512
LLM Sequence Length	128	512	512	512	512
LLM Vocab Size	256k	256k	256k	256k	256k
Codebook Size	8192	8192	8192	8192	8192

Table 7: Hyperparameters across training progression.

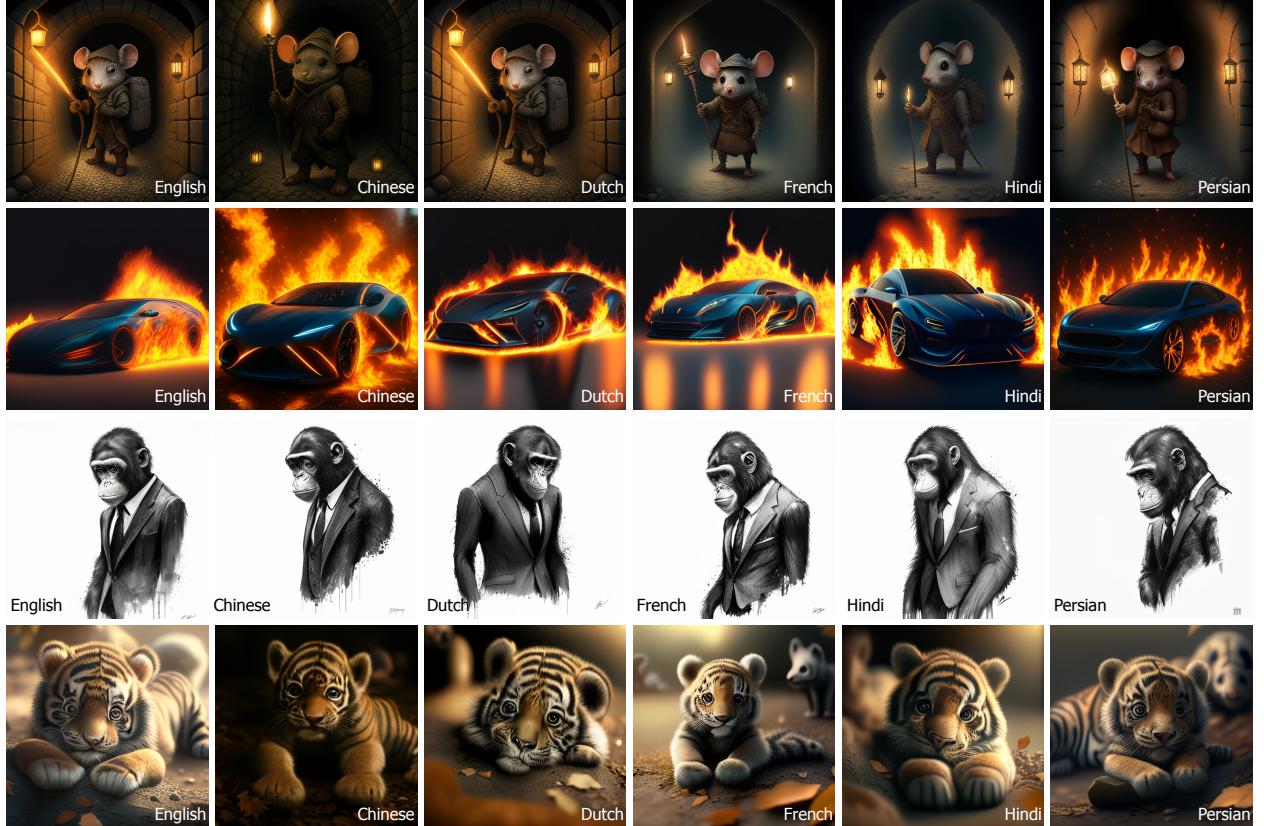


Figure 12: **Qualitative Evaluation of NEOBABEL.** Each row corresponds to a single concept expressed in six different languages: English, Chinese, Dutch, French, Hindi, and Persian. Although prompts are not shown for readability, all images were generated using translated versions of the same underlying prompt in each language. NEOBABEL consistently produces semantically aligned and visually coherent results across languages, highlighting its strong multilingual generation capabilities. We intentionally omit the prompts here due to their length, focusing instead on the visual consistency across languages.