

Aligning Constraint Generation with Design Intent in Parametric CAD

Evan Casey
Amir Khasahmadi

Tianyu Zhang
Joseph G. Lambourne

Shu Ishida
Pradeep K. Jayaraman

John R. Thompson
Karl D.D. Willis

Autodesk Research
research.autodesk.com

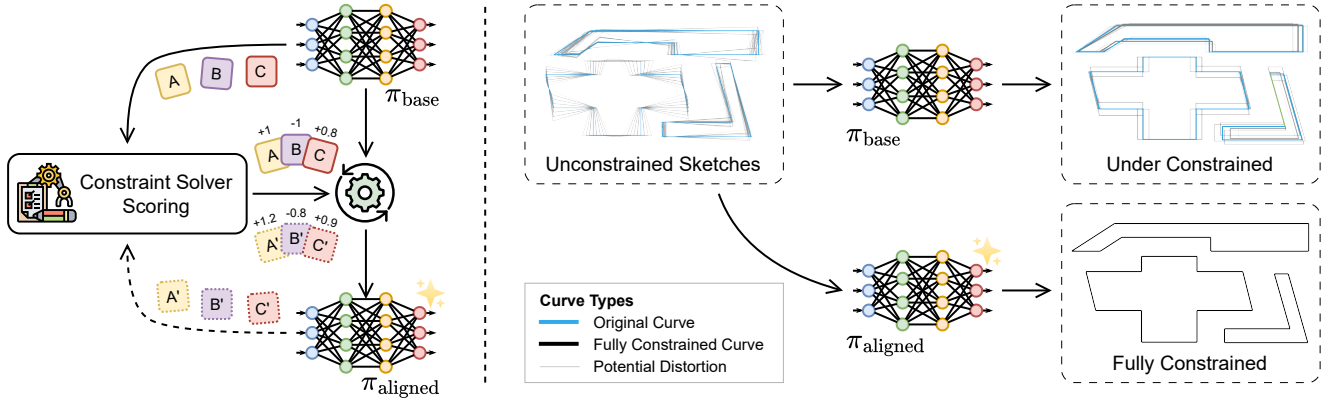


Figure 1. **Left:** a constraint solver is used to score model generated constraints $A, B, C \sim \pi_{\text{base}}$ (and $A', B', C' \sim \pi_{\text{aligned}}$). Starting with the base model π_{base} , we post-train an aligned model π_{aligned} from this feedback. **Right:** Blue lines show original primitives and gray lines show geometric distortion when dimensions vary. The aligned model π_{aligned} produces fully-constrained sketches that preserve relative geometric relationships, whereas the base model π_{base} produces under-constrained sketches that may distort the geometry in unintended ways.

Abstract

We adapt alignment techniques from reasoning LLMs to the task of generating engineering sketch constraints found in computer-aided design (CAD) models. Engineering sketches consist of geometric primitives (e.g. points, lines) connected by constraints (e.g. perpendicular, tangent) that define the relationships between them. For a design to be easily editable, the constraints must effectively capture design intent, ensuring the geometry updates predictably when parameters change. Although current approaches can generate CAD designs, an open challenge remains to align model outputs with design intent, we label this problem ‘design alignment’. A critical first step towards aligning generative CAD models is to generate constraints which fully-constrain all geometric primitives, without over-constraining or distorting sketch geometry. Using alignment techniques to train an existing constraint generation model with feedback from a constraint solver, we are able to fully-constrain 93% of sketches compared to 34% when using a naïve supervised fine-tuning (SFT) base-line and only 8.9% without alignment. Our approach can be

applied to any existing constraint generation model and sets the stage for further research bridging alignment strategies between the language and design domains. Additional results can be found in the supplementary materials.

1. Introduction

A central challenge in artificial intelligence (AI) is alignment: ensuring that AI systems produce outputs that adhere to human goals and expectations [10, 23, 27]. Although alignment of language models has been researched extensively [23, 25, 32], the application of alignment techniques to parametric design problems has yet to be studied. The use of AI in this discipline covers a broad range of areas, ranging from floor-plan layout [21, 31], to engineering design problems [8, 26, 38], to 3D generation [41]. Design problems are often visual in nature and incorporate other functional requirements, making them unique when compared with language model alignment. In this paper, we establish the problem of *design alignment* and demonstrate how this expansive problem can be made tractable by adapting

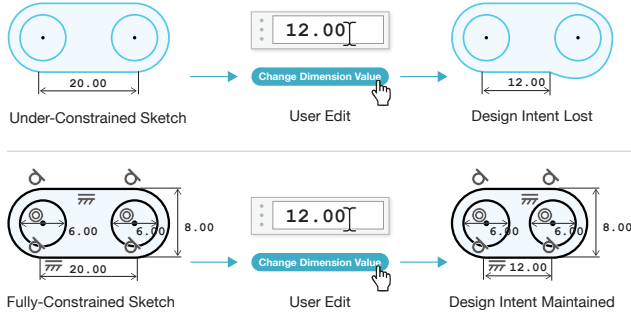


Figure 2. An illustration of design intent in CAD modeling. The bottom sketch maintains symmetry after modifying a dimension due to properly applied constraints, while the top sketch, lacking adequate constraints, becomes asymmetrical and distorted.

techniques from the alignment literature into a new context.

In language models, alignment is achieved by incorporating feedback to generate coherent, contextually appropriate responses. Similarly, with parametric CAD modeling, AI tools must be aligned with a designer’s intent by maintaining the underlying structural relationships to produce outputs that are both meaningful and functional. Otey et al. [22] define design intent as “a CAD model’s anticipated behavior when altered,” while Martin [20] characterizes it as “relationships between objects, so that a change to one propagates automatically to others.” This means that modifications of a design by an AI system should yield outcomes where the established design relationships remain intact. To that end, we define *design alignment* as the application of generative modeling alignment techniques to produce outcomes that maintain design intent.

The realization of AI systems that observe and maintain design intent has broad implications for the manufacturing and construction industries. Almost every manufactured object or structure begins as a CAD model. At the core of parametric CAD modeling are 2D engineering sketches, which can be extruded or revolved to generate 3D models. Engineering sketches are composed of geometric primitives, such as points, lines, and circles, that are organized using constraints and dimensions [7]. These constraints¹ define geometric rules, including equality, perpendicularity, and radial or linear dimensions, which collectively shape the final layout of the sketch. When applied correctly, they enable efficient modifications while preserving the original design intent.

Figure 2 illustrates the impact of constraint quality: a poorly constrained sketch loses symmetry when a dimension is changed, whereas a well-constrained sketch preserves its intended relationships. This underscores the im-

¹Throughout this paper, the term “constraints” is used in a broad sense to include both constraints (e.g., parallel) and dimensions (e.g., diameter).

portant role of constraints in maintaining design intent, and the need for AI systems that can align with this intent encoded in designs. We focus on the problem of sketch constraint generation [29] to demonstrate the adaption of alignment techniques to a design problem. Using an existing sketch constraint generation model Vitruvion [30], we align the model with algorithms that learn from feedback (Direct Preference Optimization [25], Expert Iteration [2, 33], RLOO [1], ReMax [19] and Group Relative Policy Optimization [32]) using the sketch constraint solver in Autodesk Fusion [3] as the learning signal. We optimize the models to remove all degrees of freedom in the sketches to become ‘fully-constrained’ [5], without causing sketches to be distorted, over-constrained or unsolvable. We further define these conditions in Section 3.

To the best of our knowledge, this is the first instance of alignment methods being successfully applied to a parametric CAD design task; representing an important step forward for AI-assisted design tools. We present the following contributions:

- We establish the problem of *design alignment*, in the context of generative CAD models, as a critical component of AI-assisted CAD tools. We focus on the necessary first step of alignment for engineering sketches.
- We introduce a post-training strategy for a sketch constraint generation model using feedback from a sketch constraint solver. We define novel metrics and reward functions that directly optimize a base model for improved alignment.
- We conduct extensive experiments and demonstrate alignment techniques that fully-constrain 93% of sketches compared to 34% when using a naïve supervised fine-tuning (SFT) baseline and only 8.9% without alignment. We posit that our approach is broadly applicable to other design tasks that require compilation of elements with rule-based algorithms.

2. Related Work

Engineering Sketches Engineering sketches form the 2D basis for 3D CAD models used to design mechanical parts for manufacturing. The availability of engineering sketch datasets [12, 29, 38] has enabled the development of generative models [12, 24, 30, 37] that can predict sketch geometry and/or the underlying constraints and dimensions that encode design intent. These Transformer-based [34] approaches create geometry by autoregressively generating tokens representing points and curves, then add constraints by referencing this geometry using Pointer Networks [35]. More recent approaches leverage image-based guidance [16, 39] or large language models (LLM) [15] in the constraint prediction task. Yang and Pan [42] learn to group together recurring patterns of geometric and con-

straint entities within a sketch, effectively discovering latent design concepts. However, none of these approaches explicitly optimize for preserving design intent – as a result, generated sketches may require additional manual refinement to capture the designer’s intent.

Our work builds upon these foundations by explicitly incorporating design intent as a post-training process. Instead of merely modeling the ground truth data, our method learns from constraint solver feedback, ensuring that generated sketches are geometrically plausible and structurally well-constrained. By doing so, we enable data-driven generation that aligns with design intent.

Design Alignment Beyond the language domain, alignment techniques have been used to improve and align image generation. Lee et al. [18] propose fine-tuning diffusion-based text-to-image models using human feedback, significantly improving alignment between textual prompts and generated visuals. Similarly, ImageReward [40] uses a learned reward model trained on human preference data, which guides the diffusion model fine-tuning toward images preferred by human evaluators. Extending this idea, Black et al. [6] reframes image generation as a sequential RL task, introducing Denoising Diffusion Policy Optimization (DDPO) to optimize complex user-defined objectives directly for alignment without explicit human annotation. Within the design domain, few works applied generative modeling alignment techniques to produce outcomes that maintain design intent. GearFormer [11] leveraged differentiable sampling to enforce preferences in the solutions for mechanical configuration design problems, however, this approach does not work with rewards that require black-box solvers in the loop. In concurrent work, e-SimFT [9] leveraged preference data obtained from physics simulations to enhance the exploration of the Pareto front in a multi-preference setting, improving the optimality of solutions generated with GearFormer. We believe the problem of *design alignment* will become a critical part of future generative models aimed at design applications.

Fine-tuning LLMs with RL Reinforcement Learning from Human Feedback (RLHF) has emerged as a cornerstone approach for aligning large language models (LLMs) with human preferences. In RLHF [4, 10, 23, 28], a learned reward model is usually trained to capture human preferences and provides the learning signal that the policy is trained on. Other methods such as Direct Preference Optimization (DPO) [25] and Reinforced Self-Training (ReST) [33] use a simpler approach of optimizing the model directly from model generated data without the use of a learned reward model. While DPO learns from ranked pairs of generated model outputs from human annotators, ReST uses rejection sampling to remove incorrect generations and trains the model standard cross-entropy loss

on the correct samples. In this paper, we refer to the approach of using fine-tuning on high return responses as Expert Iteration (ExIt). We broadly refer to algorithms that learn from ranked/filtered model outputs (such as DPO and Expert Iteration) as Preference Optimization (PO).

More recently, a large body of work has focused on the task of teaching LLMs to solve reasoning tasks with reinforcement learning [13, 14, 17, 33], such as math and coding, which can be checked with rule-based systems. Specifically we are inspired by approaches which forego the use of a learned reward model and directly learn from verifiable rewards. We build off of several approaches that have shown success when applied to reasoning LLMs – these include Group Relative Policy Optimization (GRPO) [32], ReMax [19], and Reinforce Leave-One-Out (RLOO) [1]. Both GRPO and RLOO estimate the baseline via the average reward of multiple sampled outputs instead of learned value model but differ in how they apply the KL divergence penalty, advantage normalization, and PPO-style reward clipping. ReMax [19] also obviates the need for a learned value model but instead estimates the baseline from the argmax result (greedy sampling).

In this paper, we adapt the aforementioned post-training methods—DPO, Expert Iteration, RLOO, GRPO, ReMax, for use in aligning a constraint generation model using feedback from a constraint solver. In Section 4 we provide additional details on the ranking/filtering criteria for the Preference Optimization (PO) algorithms and the reward design for the RL algorithms (GRPO, RLOO, ReMax).

3. Problem

Sketch constraining is a fundamental component of parametric CAD modeling, where geometric relationships define the structure and behavior of the sketch. Applying constraints ensures stability and editability, allowing for parametric modifications that align with the design intent. Automating constraint generation requires producing a valid and efficient set of constraints that fully define a given sketch while avoiding unnecessary redundancy or conflicts.

The sketch constraining problem can be formulated as a sequence modeling task similar to natural language generation, where constraints are predicted autoregressively. Given the sketch geometry as input, the model generates a sequence of constraints in the order they will be applied. Tokens in the sequence represent either a constraint (e.g., coincident, parallel, perpendicular), a dimension (e.g., horizontal, vertical, radial), or a pointer to one of the input geometric entities [35].

In parametric CAD, sketch geometry is modified using a constraint solver. The updated sketch respects any present constraints while moving the geometry to reflect changes to the dimension parameters. Unlike natural language, which is inherently sequential and follows flexible grammar rules,

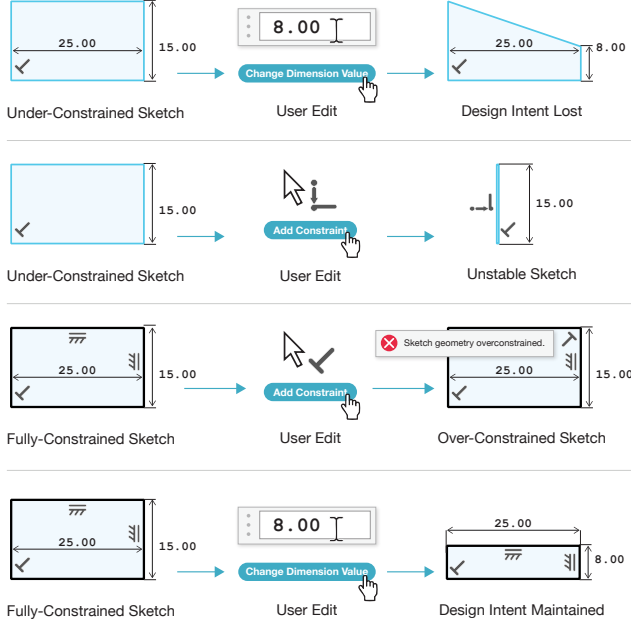


Figure 3. Comparing different outcomes when a designer updates a sketch parameter on the constrained sketch. **First row:** an under-constrained sketch only preserves a subset of the geometric relationships. **Second row:** the sketch is unstable, adding a coincident constraint flattens the geometry of the sketch. **Third row:** the sketch is over-constrained, causing the sketch to be uneditable. **Fourth row:** all geometric relationships are maintained after the parameter is updated.

sketch constraints must adhere to strict geometric principles to ensure structural validity. A set of constraints may be incorrect for a variety of reasons: they can reference the wrong primitives for the constraint type, be redundant, specify inconsistent geometric relationships, or cause unexpected geometric distortions.

We formally define five conditions describing the state of a sketch after applying constraints. These conditions are not mutually exclusive. A sketch may satisfy one or more conditions depending on the applied constraints. Figure 3 provides illustrative examples of each condition.

Under-constrained (UC) A sketch containing primitives retains some unconstrained degrees of freedom, resulting in incomplete specification of their positions or dimensions.

Fully-constrained (FC) A sketch in which all primitives have their degrees of freedom completely determined, removing positional or dimensional ambiguities.

Over-constrained (OC) A sketch primitives have more constraints applied than degrees of freedom, potentially leading to conflicts. Note some over-constrained sketches remain solvable if the constraints are consistent and do not conflict with each other [7].

Not solvable A sketch that cannot achieve a valid solution due to contradictory or redundant constraints, leading to an impossible or conflicting geometry.

Stability We discretize the sketch plane into a grid and classify a sketch as **unstable** if the positions of primitives after constraint solving shift into different cells. The number of cells (bins) on each axis determines the sensitivity of this measurement.

Our objective is to align the model toward generating constraint sets that yield fully-constrained sketches while minimizing cases of under-constrained, over-constrained, not solvable, or instability. This is a necessary prerequisite toward the ultimate goal of generating constraints that preserve the original sketch design intent when the designer varies the geometric parameters.

4. Method

In this section, we outline the post-training techniques used for aligning a constraint generation model with feedback from a constraint solver. The approaches are grouped into three categories: supervised learning methods, preference-based optimization methods, and RL methods. Figure 4 illustrates the high-level workflow of this work.

4.1. Constraint Solver

To evaluate model-generated constraints, we integrate the commercially available constraint solver in Autodesk Fusion. This provides industrial-grade accuracy and robustness, allowing us to assess sketch conditions defined in Section 3, whether applying a given set of constraints yields a stable, fully-constrained sketch or leads to over-constraint and unsolvability. Additionally, the Fusion constraint solver adjusts geometry to resolve all constraints, enabling alignment checks with design intent. On average the constraint solver takes 0.1-0.2 seconds to evaluate a sketch, although difficult sketches can be in the tens of seconds. For computational efficiency we automatically deem the sketch unsolvable if the sketch takes longer than two seconds to solve.

4.2. Supervised Learning

We pre-train Vitruvion [30] as our base model using the same procedure described in their paper with a next-token prediction objective. Given an input sequence of geometric primitives, the model is trained to predict the next correct constraint or dimension based on the ground truth data. Further details about our implementation of the Vitruvion architecture and training procedure are provided in the appendix.

Since the majority of the ground truth data contains under-constrained or over-constrained sketches, we additionally perform supervised fine-tuning (SFT). During SFT, the training data is limited to sketches verified by the constraint solver as solvable, fully-constrained, stable, and free

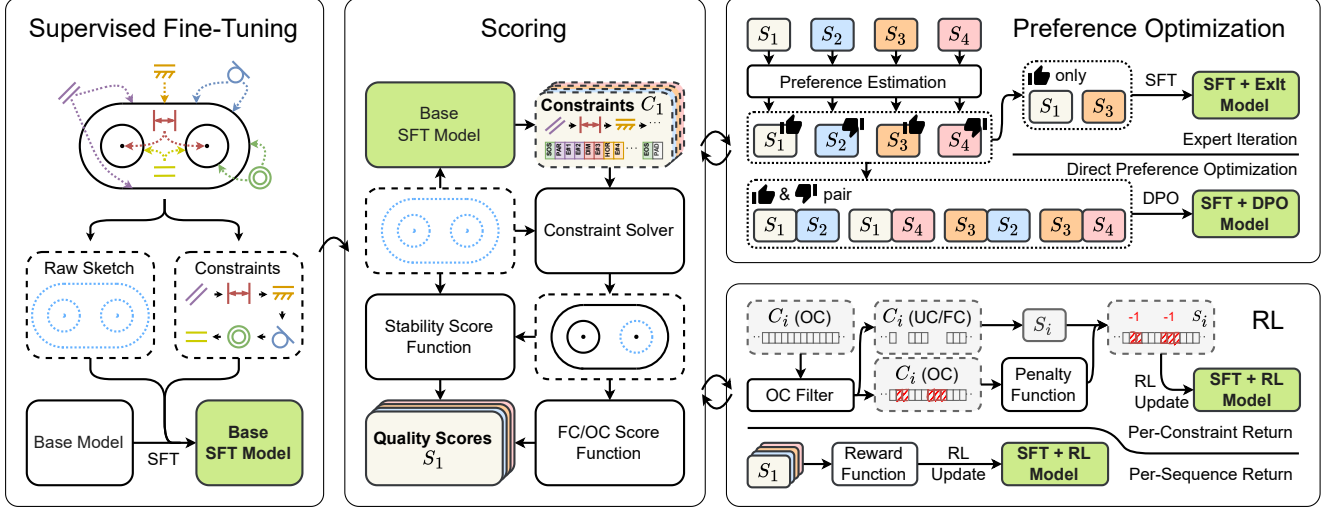


Figure 4. Illustration of the proposed alignment workflow for constraint generation models. **Left:** A base constraint-generation model is first fine-tuned using supervised fine-tuning (SFT). **Middle:** Generated constraint sequences C_i are evaluated using a constraint solver, which provides feedback on sketch stability and fully- and over-constrained statuses, forming the quality score vector S_i . **Right:** Two groups of alignment methods leveraging solver feedback: preference-based optimization (PO) and reinforcement learning (RL). PO uses S_i to construct training data, iteratively improving constraint prediction quality. RL methods assign per-constraint and per-sequence rewards to C_i based on solver feedback and S_i to incentivize the generation of constraints, making sketches stable and fully-constrained.

of over-constrained conditions, ensuring the model explicitly learns from ideal examples.

4.3. Preference-Based Optimization (PO)

Expert Iteration (ExIt) alternates between expert improvement and policy distillation. Following [14, 33], we use temperature sampling combined with rejection sampling to generate high-quality candidate constraint sequences. Specifically, during the exploration step, we initialize the policy model $\pi_{\theta_{t=0}}$ from the SFT model, and for each sketch query q in the training set, we sample $K = 8$ candidate constraint sequences τ at temperature $T = 1.0$. The training dataset is constructed by discarding sequences that are under-constrained, over-constrained, or unsolvable solutions. This process is repeated $N = 2$ times over the dataset, and the policy is trained using cross-entropy loss:

$$\mathbb{E}_{(q, \tau) \sim \mathcal{D}} [\log \pi_{\theta_t}(\tau|q)] \quad (1)$$

where π_{θ_t} is updated after the distillation phase of every training round.

Direct Preference Optimization (DPO) learns from pairwise preference data, approximating an implicit reward via a reparameterized Bradley-Terry model [25]. Similar to ExIt, we construct the training dataset by sampling $K = 8$ constraint sequence completions τ for each sketch query q from the policy model π_{θ_t} at temperature $T = 1.0$. Pairs (τ_w, τ_l) are ranked based on the differences in the percentage of fully-constrained curves between τ_w and τ_l . The op-

timization objective is:

$$\mathbb{E}_{(q, \tau_w, \tau_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta_t}(\tau_w|q)}{\pi_{\theta_r}(\tau_w|q)} - \beta \log \frac{\pi_{\theta_t}(\tau_l|q)}{\pi_{\theta_r}(\tau_l|q)} \right) \right] \quad (2)$$

where β is a hyperparameter, and σ is the logistic function. This formulation ensures the learned policy π_{θ} aligns with the ranking preference of fully-constrained sketches while maintaining proximity to the reference policy π_{θ_r} .

We initialize $\pi_{\theta_{t=0}}$ from the SFT model and repeat the entire process $N = 2$ times, updating π_{θ_t} with the new policy after every training iteration. Additional hyperparameters and ranking details can be found in the appendix.

4.4. Reinforcement Learning (RL)

4.4.1. Reward design

Unlike natural language tasks, where human preferences are ambiguous and ill-defined, the stability and solvability of sketch constraints can be verified, making it compatible with RL methods that directly optimize for mechanically defined rewards without a learned preference model. We define the rewards used for RL as follows:

- Rewards for valid constraint sequence τ :
 - $r_{\text{curves}}(\tau)$: % of fully-constrained curves over all curves,
 - $r_{\text{points}}(\tau)$: % of fully-constrained points over all points,
 - r_{unstable} : penalty for unstable sketches,
- Rewards for invalid constraint sequence:
 - r_{NS} : penalty for not solvable sketches,

r_{OC} : penalty for over-constrained sketches,

r_{F} : penalty for sketches resulting in other failures.

The overall sequence-wise reward $R(\tau)$ is the sum of $r_{\text{curves}}(\tau)$, $r_{\text{points}}(\tau)$, and conditionally r_{unstable} for valid sketches, and either r_{NS} , r_{OC} or r_{F} for invalid sketches according to the failure mode.

We additionally define a constraint-wise penalty to provide granular feedback on cases where the constraint sequence causes sketches to be over-constrained or fully-constrained. A constraint solver iteratively attempts to add each generated constraint one-by-one, dropping any problematic constraints that caused the sketch to be over-constrained or not-solvable. In training, we add a constant of -1 loss penalty directly to the per-token log likelihood loss for the problematic constraints.

4.4.2. RL fine-tuning formulation

We formulate constraint generation fine-tuning as follows; given a dataset of sketch queries $D = \{q_i\}_{i=1}^N$ and reward function $R(\tau)$ using the constraint solver and reward design in Section 4.4.1, learn a policy $\pi_\theta(\tau|q)$ that generates a sequence of constraints τ for sketch q , such that it maximizes the expected rewards $\mathbb{E}_{q_i \sim D, \tau_i \sim \pi_\theta(\cdot|q_i)}[R(\tau_i)]$.

4.4.3. Policy gradient methods

For RLHF which uses a pre-trained policy, not all the complexity of RL algorithms is necessary. This allows the algorithms to be simplified and the number of learnable components to be reduced, contributing to performance improvement. We considered three policy gradient algorithms: ReMax [19], RLOO [1], and GRPO [32]. Unlike PPO [28], which treats each token generation as an action, these algorithms treat generation of a sequence as a single action and adopt a REINFORCE with baselines approach [36]. We apply these to optimize the constraint generation policy.

ReMax [19] uses the rewards corresponding to sequences generated by a greedy (argmax) policy as a baseline to normalize the rewards of sequences sampled from the policy.

Considering sequence τ sampled from policy $\pi_\theta(\tau|q)$, sequence $\tau^* = \text{argmax}_\tau \pi_\theta(\tau|q)$ greedily sampled by taking an argmax of the policy, and corresponding rewards r and r^* , the policy gradient objective for ReMax is:

$$\mathbb{E}_{\tau \sim \pi} [(r - r^*) \nabla \log \pi_\theta(\tau|q)]. \quad (3)$$

REINFORCE-Leave-One-Out (RLOO) [1] samples G number of constraint sequences for every sketch query. The baseline for each sample is evaluated as the mean of the rewards for all other samples in the group.

For G number of sequences $\{\tau_g\}_{g=1}^G$ sampled from policy $\pi_\theta(\tau|q)$ for a given sketch query q , the policy gradient objective of RLOO is:

$$\mathbb{E}_{\{\tau\} \sim \pi} \left[\frac{1}{G} \sum_{g=1}^G \left[\left(r_g - \text{mean}(\{r_i\}_{i \neq g}^G) \right) \nabla \log \pi_\theta(\tau_g|q_g) \right] \right]. \quad (4)$$

Group Relative Policy Optimization (GRPO) [32] uses group-based baseline estimation, like RLOO, but the mean is taken over all reward samples in the group. It uses a clipped policy optimization objective similarly to PPO [28], as well as a low-variance KL regularization term.

For G number of sequences $\{\tau_g\}_{g=1}^G$ sampled from the reference policy $\pi_{\theta_t}(\tau|q)$ for a given sketch query q , letting $\rho_g = \frac{\nabla \pi_\theta(\tau_g|q)}{\pi_{\theta_t}(\tau_g|q)}$, the optimization objective of GRPO is:

$$\mathbb{E}_{\{\tau_g\} \sim \pi} \left[\min(\rho_g A_g, \text{clip}(\rho_g, 1 - \epsilon, 1 + \epsilon) A_g) - \beta \mathbb{D}_{\text{KL}}(\pi_\theta || \pi_{\theta_t}) \right],$$

$$\mathbb{D}_{\text{KL}}(\pi_\theta || \pi_{\theta_t}) = \frac{1}{\rho_g} + \log \rho_g - 1, \quad A_g = \frac{r_g - \text{mean}(\{r_g\}_{g=1}^G)}{\text{std}(\{r_g\}_{g=1}^G)}, \quad (5)$$

where ϵ and β are hyper-parameters for the clipped policy optimization and KL regularization terms, respectively.

For ReMax and RLOO, we also added a small KL penalty term to the rewards to discourage divergence from the reference policy. GRPO applies group-normalization on the advantage, which we also applied in RLOO. For ReMax, we batch-normalized the advantage. For all algorithms, π_θ is initialized from the SFT model. Further implementation details of the algorithms can be found in the appendix.

5. Experiments

5.1. Dataset

We train our models on SketchGraphs [29], a large-scale dataset of CAD sketches created in Onshape. SketchGraphs captures real-world parametric modeling workflows, providing geometry and constraint construction operations from actual design steps. However, the dataset was not originally designed for direct constraint inference; only 8.27% of its sketches are fully-constrained, making it imperfect for training the constraint generation model.

For computational feasibility, we deduplicate and filter sketches, retaining only those with at most 16 geometric primitives and 64 constraints, yielding a dataset of 2.8 million unique sketches. Certain constraint types, such as symmetry, are excluded to simplify the learning task, ensuring the focus remains on constraints most relevant to engineering design. We also convert Onshape sketches into the Fusion sketch format to utilize the solver in Fusion. Additional details are provided in the appendix.

Another challenge lies in how SketchGraphs positions primitives. Rather than being placed in valid, constraint-satisfying layouts, primitives often have arbitrary coordinates that do not reflect a solved state. We therefore prepro-

Table 1. Sketch constraint generation results for Fully Constrained (FC), Under Constrained (UC), Over Constrained (OC), not solvable, and stability for the base model, SFT model, and aligned models. Results are computed over 8 samples per sketch with a temperature of 1.0. Numbers following \pm indicate the standard deviation.

Model	% FC \uparrow	% UC \downarrow	% OC \downarrow	% Not solvable \downarrow	% Stable (bins=4) \uparrow
Vitruvion (base)	8.87 ± 0.09	71.38 ± 0.20	16.83 ± 0.17	3.05 ± 0.03	92.15 ± 0.06
SFT	34.24 ± 0.09	46.61 ± 0.15	15.30 ± 0.08	3.85 ± 0.03	92.48 ± 0.05
Iterative DPO	64.91 ± 0.11	14.97 ± 0.13	12.47 ± 0.09	7.64 ± 0.07	87.63 ± 0.09
Expert Iteration	71.70 ± 0.13	13.38 ± 0.13	7.25 ± 0.06	7.67 ± 0.07	85.50 ± 0.10
ReMax	79.84 ± 0.09	15.86 ± 0.07	1.49 ± 0.01	2.82 ± 0.02	75.77 ± 0.05
RLOO	93.05 ± 0.03	3.55 ± 0.02	2.15 ± 0.01	1.25 ± 0.01	89.16 ± 0.02
GRPO	91.59 ± 0.03	4.18 ± 0.03	1.94 ± 0.01	2.28 ± 0.02	88.28 ± 0.03

cess the dataset using Fusion to resolve each sketch’s primitives according to its constraints, ensuring that geometry and constraints match before training.

5.2. Quantitative Results

In Table 1, we list results comparing the performance of each alignment method with respect to the five sketch conditions described in Section 3. The base model is able to fully-constrain sketches only 8.87% of the time, consistent with the dataset distribution where only 8.27% of sketches are fully constrained. We find that RLOO and GRPO perform similarly, giving the best performance at fully-constraining sketches 93.05% and 91.59% of the time, respectively. They have the lowest indicents of over-constrained or unsolvable results and maintain stability rates that are on par with other methods.

Iterative DPO and ExIt significantly improve upon the base and SFT models but still fall short of the performance achieved by policy gradient-based RL methods. We attribute this gap to the online nature of policy gradient-based RL, which continuously refines the policy through feedback while actively exploring a broader range of solutions. In contrast, Iterative DPO and ExIt are offline methods and rely on predefined ranking and filtering signals to generate training data, which limits their ability to explore the so-

lution space. The superior performance of online RL underscores its advantage in directly optimizing the shaped rewards from the constraint solver.

Table 2 demonstrates the evaluation result of our methods in a few-shot inference setting, where the model has K attempts to generate a valid solution. Since the constraint solver can verify whether a solution is correct, this scenario reflects real-world use cases where the goal is to maximize the likelihood of producing an acceptable solution within a fixed inference budget of K samples (Pass@ K). We find that while the policy gradient-based RL methods still have the overall highest performance, increasing the number of samples K has a comparatively small impact on performance compared to the other methods.

5.3. Qualitative Results

We present qualitative comparisons of sketches generated by different alignment methods in Figure 5. Curves are colored black when constrained and blue when not. Methods leveraging constraint solver feedback (Iterative DPO, ExIt, ReMax, RLOO, GROO) demonstrate a promising trend towards generating fully-constrained sketches, whereas the Base and SFT models typically leave sketches under-constrained. However, across different alignment algorithms, we observe substantial variance in the degree of geometric distortion introduced by the aligned models.

Specifically, columns A and B depict simple sketches mainly composed of horizontal and vertical lines, for which all solver-feedback methods consistently yield fully-constrained, stable results. In contrast, Column C presents a challenging sketch due to the absence of appropriate constraints for oblique lines. Achieving stability in this scenario typically requires many dimensions but few constraints; however, models are optimized to generate more constraints to align with parametric CAD design principles.

Sketches in columns D through H include arcs. Column E is particularly notable as all solver-feedback methods achieve a fully-constrained condition, yet only RLOO produces visually stable results. Further examination reveals that distortions induced by other models were still classi-

Table 2. Sketch constraint generation results for Pass@1 and Pass@8 across the post-training algorithms. We define a successful result as fully-constrained, not over-constrained, solvable, and stable at 4 bins. Results are generated with temperature of 1.0.

Model	Pass@1	Pass@8
Vitruvion (base)	8.53	20.47
SFT	33.32	42.62
Iterative DPO	59.38	68.32
Expert Iteration	64.09	72.56
ReMax	62.74	65.89
RLOO	83.57	84.96
GRPO	81.49	83.42

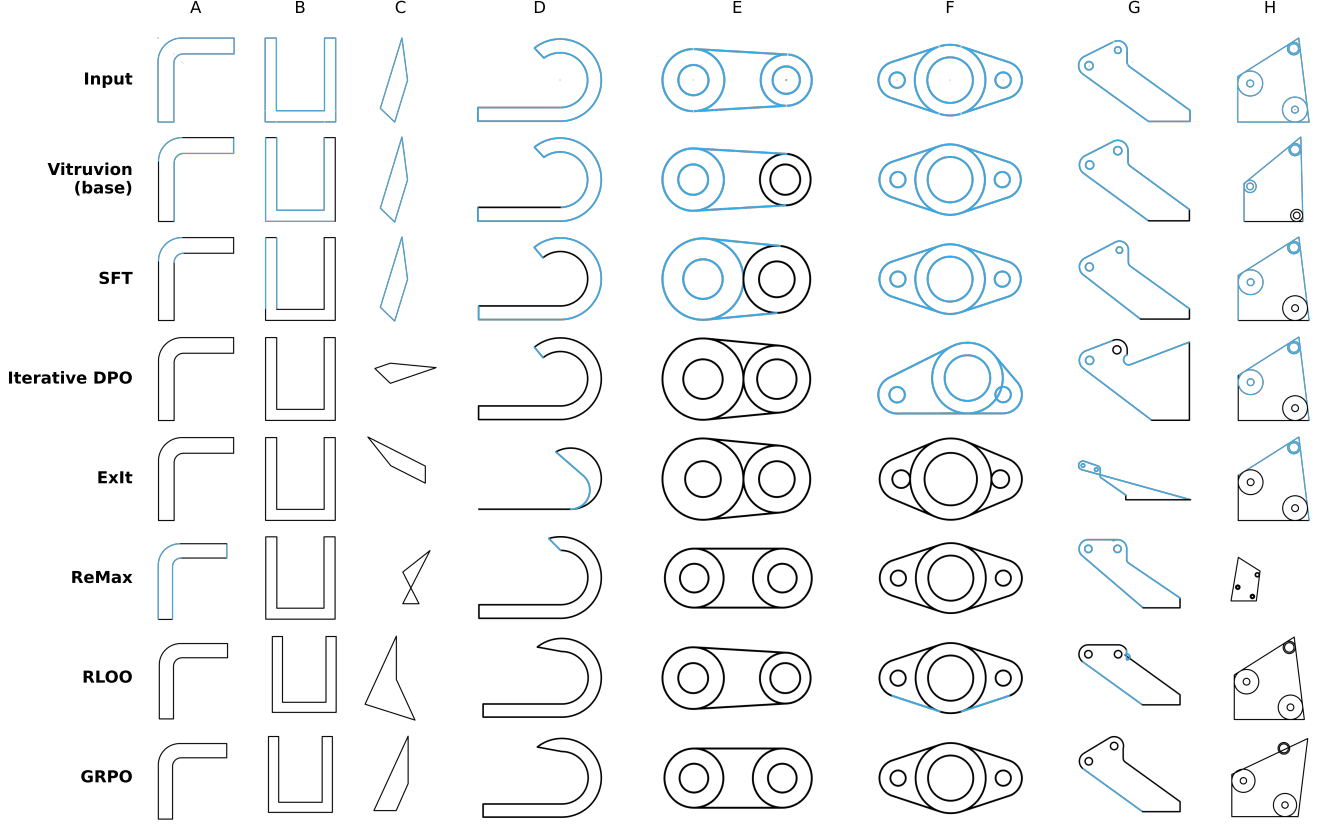


Figure 5. Visual comparison of sketches with constraints generated using a baseline model and with different post-training approaches applied. Curves are colored black when constrained and blue when not.

fied as stable due to our bin size ($bins = 4$). Column G presents unique challenges with an isolated point, causing confusion in the model and demanding extensive use of non-horizontal/vertical constraints to make the sketch fully-constrained and stable.

In summary, our results indicate that models more easily fully-constrain sketches consisting primarily of horizontal and vertical lines without distortion. In contrast, sketches involving oblique lines increase the challenge of maintaining stability, and arcs further complicate achieving fully-constrained status. Among the tested methods, RLOO and GRPO exhibit the strongest overall performance in constraint satisfaction and geometric stability. Additional qualitative results can be found in the supplementary materials.

6. Limitations

We now outline several limitations of our work suitable for future research. First, our post-training approach has not been optimized for speed. Currently it takes ~ 3 days to train a single epoch with the SketchGraphs dataset using a 8xH100 GPU configuration. This is primarily due to the frequent interactions with the CPU-based constraint solver and

the fact that solve times can be highly varied. Roughly half of the training time is spent on GPU computation and half on detokenization and solver interaction. We expect custom optimizations could significantly reduce training time. Second, our experiments focus on sketches of moderate complexity, with a maximum of 16 geometric primitives and 64 constraints. Evaluating performance on larger, more complex sketches requires further exploration. Finally, our alignment approach currently uses feedback signals from a constraint solver to represent design intent. Incorporating subjective human preferences and explicit design intent into alignment objectives remains a promising area for future work.

7. Conclusion

We demonstrated the adaption of alignment strategies from language modeling to preserve design intent in parametric CAD sketches. By using feedback from a constraint solver as a learning signal, we show the feasibility and value of alignment in parametric CAD tasks. In doing so, we pave the way for future AI-assisted design tools that incorporate *design alignment*.

References

- [1] Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms. *arXiv preprint arXiv:2402.14740*, 2024. 2, 3, 6
- [2] Thomas W. Anthony, Zheng Tian, and David Barber. Thinking fast and slow with deep learning and tree search. In *Neural Information Processing Systems*, 2017. 2
- [3] Autodesk. *Sketches in Fusion*, 2014. 2
- [4] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, John Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022. 3
- [5] Bernhard Bettig and Christoph M. Hoffmann. Geometric constraint solving in parametric computer-aided design. *Journal of Computing and Information Science in Engineering*, 11(2):021001, 2011. 2
- [6] Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023. 3
- [7] William Bouma, Ioannis Fudos, Christoph Hoffmann, Jiazhen Cai, and Robert Paige. Geometric constraint solver. *Computer-Aided Design*, 27(6):487–501, 1995. 2, 4
- [8] Wei Chen, Kevin Chiu, and Mark D Fuge. Airfoil design parameterization and optimization using bézier generative adversarial networks. *AIAA journal*, 58(11):4723–4735, 2020. 1
- [9] Hyunmin Cheong, Mohammadmehdi Ataei, Amir Hosein Khasahmadi, and Pradeep Kumar Jayaraman. e-simft: Alignment of generative models with simulation feedback for pareto-front design exploration. *arXiv preprint arXiv:2502.02628*, 2025. 3
- [10] Paul F. Christiano, Jan Leike, Tom B. Brown, et al. Deep reinforcement learning from human preferences. In *NeurIPS*, 2017. 1, 3
- [11] Yasaman Etesam, Hyunmin Cheong, Mohammadmehdi Ataei, and Pradeep Kumar Jayaraman. Deep generative model for mechanical system configuration design. *arXiv preprint arXiv:2409.06016*, 2024. 3
- [12] Yaroslav Ganin, Sergey Bartunov, Yujia Li, Ethan Keller, and Stefano Saliceti. Computer-aided design as language. *NeurIPS*, 34:5885–5897, 2021. 2
- [13] Jonas Gehring, Kunhao Zheng, Jade Copet, Vegard Mella, Taco Cohen, and Gabriel Synnaeve. Rlf: Grounding code llms in execution feedback with reinforcement learning. *arXiv preprint arXiv:2410.02089*, 2024. 3
- [14] Alex Havrilla, Yuqing Du, Sharath Chandra Raparthy, Christoforos Nalmpantis, Jane Dwivedi-Yu, Maksym Zhuravinskyi, Eric Hambro, Sainbayar Sukhbaatar, and Roberta Raileanu. Teaching large language models to reason with reinforcement learning. *arXiv preprint arXiv:2403.04642*, 2024. 3, 5
- [15] Benjamin T Jones, Felix Hähnlein, Zihan Zhang, Maaz Ahmad, Vladimir Kim, and Adriana Schulz. A solver-aided hierarchical language for llm-driven cad design. *arXiv preprint arXiv:2502.09819*, 2025. 2
- [16] Ahmet Serdar Karadeniz, Dimitrios Mallis, Nesryne Mejri, Kseniya Cherenkova, Anis Kacem, and Djamila Aouada. Davinci: A single-stage architecture for constrained cad sketch inference. *arXiv preprint arXiv:2410.22857*, 2024. 2
- [17] Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. T\”ulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024. 3
- [18] Kimin Lee, Hao Liu, Moonkyung Ryu, Olivia Watkins, Yuqing Du, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, and Shixiang Shane Gu. Aligning text-to-image models using human feedback. *arXiv preprint arXiv:2302.12192*, 2023. 3
- [19] Ziniu Li, Tian Xu, Yushun Zhang, Zhihang Lin, Yang Yu, Ruoyu Sun, and Zhi-Quan Luo. Remax: A simple, effective, and efficient reinforcement learning method for aligning large language models. In *ICML*, 2024. 2, 3, 6
- [20] Dave Martin. What is design intent?, 2023. Accessed: January 19, 2023. 2
- [21] Nelson Nauata, Kai-Hung Chang, Chin-Yi Cheng, Greg Mori, and Yasutaka Furukawa. House-gan: Relational generative adversarial networks for graph-constrained house layout generation. In *ECCV*, pages 162–177. Springer, 2020. 1
- [22] Jeffrey Otey, Pedro Company, Manuel Contero, and Jorge D. Camba. Revisiting the design intent concept in the context of mechanical cad education. *Computer-Aided Design and Applications*, 15(1):47–60, 2018. 2
- [23] Long Ouyang, Jeff Wu, Xu Jiang, et al. Training language models to follow instructions with human feedback. In *NeurIPS*, 2022. 1, 3
- [24] Wamiq Para, Shariq Bhat, Paul Guerrero, Tom Kelly, Niloy Mitra, Leonidas J Guibas, and Peter Wonka. Sketchgen: Generating constrained cad sketches. *NeurIPS*, 34:5077–5088, 2021. 2
- [25] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *NeurIPS*, 36:53728–53741, 2023. 1, 2, 3, 5
- [26] Lyle Regenwetter, Amin Heyrani Nobari, and Faez Ahmed. Deep generative models in engineering design: A review. *Journal of Mechanical Design*, 144(7):071704, 2022. 1
- [27] Stuart Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach*. Prentice Hall, 3rd edition, 2010. 1
- [28] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. 3, 6
- [29] Ari Seff, Yaniv Ovadia, Wenda Zhou, and Ryan P. Adams. SketchGraphs: A large-scale dataset for modeling relational geometry in computer-aided design. In *ICML 2020 Workshop on Object-Oriented Learning*, 2020. 2, 6
- [30] Ari Seff, Wenda Zhou, Nick Richardson, and Ryan P Adams. Vitruvion: A generative model of parametric cad sketches. In *ICLR*, 2021. 2, 4

- [31] Mohammad Amin Shabani, Sepidehsadat Hosseini, and Yasutaka Furukawa. Housediffusion: Vector floorplan generation via a diffusion model with discrete and continuous denoising. In *CVPR*, pages 5466–5475, 2023. [1](#)
- [32] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. DeepSeekMath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024. [1](#), [2](#), [3](#), [6](#)
- [33] Avi Singh, John D Co-Reyes, Rishabh Agarwal, Ankesh Anand, Piyush Patil, Xavier Garcia, Peter J Liu, James Harrison, Jaehoon Lee, Kelvin Xu, Aaron T Parisi, Abhishek Kumar, Alexander A Alemi, Alex Rizkowsky, Azade Nova, Ben Adlam, Bernd Bohnet, Gamaleldin Fathy Elsayed, Hanie Sedghi, Igor Mordatch, Isabelle Simpson, Izzeddin Gur, Jasper Snoek, Jeffrey Pennington, Jiri Hron, Kathleen Kenealy, Kevin Swersky, Kshiteej Mahajan, Laura A Culp, Lechao Xiao, Maxwell Bileschi, Noah Constant, Roman Novak, Rosanne Liu, Tris Warkentin, Yamini Bansal, Ethan Dyer, Behnam Neyshabur, Jascha Sohl-Dickstein, and Noah Fiedel. Beyond human data: Scaling self-training for problem-solving with language models. *Transactions on Machine Learning Research*, 2024. [2](#), [3](#), [5](#)
- [34] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 30, 2017. [2](#)
- [35] Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. Pointer networks. In *Proceedings of the 29th International Conference on Neural Information Processing Systems - Volume 2*, page 2692–2700, Cambridge, MA, USA, 2015. MIT Press. [2](#), [3](#)
- [36] Lex Weaver and Nigel Tao. The optimal reward baseline for gradient-based reinforcement learning. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, page 538–545, San Francisco, CA, USA, 2001. Morgan Kaufmann Publishers Inc. [6](#)
- [37] Karl DD Willis, Pradeep Kumar Jayaraman, Joseph G Lambourne, Hang Chu, and Yewen Pu. Engineering sketch generation for computer-aided design. In *CVPRW*, pages 2105–2114, 2021. [2](#)
- [38] Karl D. D. Willis, Yewen Pu, Jieliang Luo, Hang Chu, Tao Du, Joseph G. Lambourne, Armando Solar-Lezama, and Wojciech Matusik. Fusion 360 gallery: A dataset and environment for programmatic cad construction from human design sequences. *ACM TOG*, 40(4), 2021. [1](#), [2](#)
- [39] Sifan Wu, Amir Hosein Khasahmadi, Mor Katz, Pradeep Kumar Jayaraman, Yewen Pu, Karl Willis, and Bang Liu. Cad-vm: Bridging language and vision in the generation of parametric cad sketches. In *ECCV*, pages 368–384. Springer, 2024. [2](#)
- [40] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36:15903–15935, 2023. [3](#)
- [41] Xiang Xu, Pradeep Kumar Jayaraman, Joseph G Lambourne, Karl DD Willis, and Yasutaka Furukawa. Hierarchical neural coding for controllable cad model generation. In *ICML*, pages 38443–38461, 2023. [1](#)
- [42] Yuezhi Yang and Hao Pan. Discovering design concepts for cad sketches. *NeurIPS*, 35:28803–28814, 2022. [2](#)

Appendix

A. Parametric CAD Sketches

We discuss additional details related to our dataset of parametric CAD sketches and constraint solver.

A.1. Background

Parametric CAD fundamentally relies on sketches as the basis for generating complex 3D geometries. Sketches are formed from geometric primitives such as points, lines, arcs, and circles. By imposing constraints (e.g., tangency, perpendicularity, parallelism) and dimensions (e.g., linear, angular, radial), these primitives become systematically interlinked, preserving design intent through iterative modifications. A dedicated constraint solver manages this network of relationships, using numerical methods to maintain consistency and automatically adjust dependent elements when any single parameter changes.

In Figure A.1, tangent constraints (blue) reference a line and an arc, horizontal constraints (orange) reference two lines, and linear dimensions (red) reference two points. Such definitions encode both geometric relationships and key measurements, allowing the solver to propagate updates throughout the model. This approach reduces the need for manual rework by ensuring that changing one dimension, such as the distance between two points or the radius of an arc, will automatically update the entire sketch. This allows designers to iterate rapidly while maintaining the design intent embedded in the sketch.

A.2. Fusion Sketch Representation

The Fusion format organizes sketch elements into a hierarchical, structured representation, wherein a sketch is defined by a set of parametric geometric primitives and a set of explicit constraints between those primitives. Each geometric primitive (line, arc, circle, point, etc.) is described by its intrinsic parameters (e.g., endpoint coordinates for a line, center and radius for a circle). Alongside the primitives, the sketch includes constraints (e.g., coincident points, perpendicular or parallel lines) that impose geometric relationships to be satisfied simultaneously. These constraints serve to preserve design intent: for instance, a coincidence constraint can lock the endpoint of a line onto a circle’s circumference, or an equal-length constraint can enforce that two segments remain the same length.

Structuring the sketch with primitives and constraints yields a rich, relational format rather than a flat drawing. The representation can be viewed as a bipartite graph, where primitive nodes carry geometric parameters and constraint edges specify relationships linking one or more primitives.

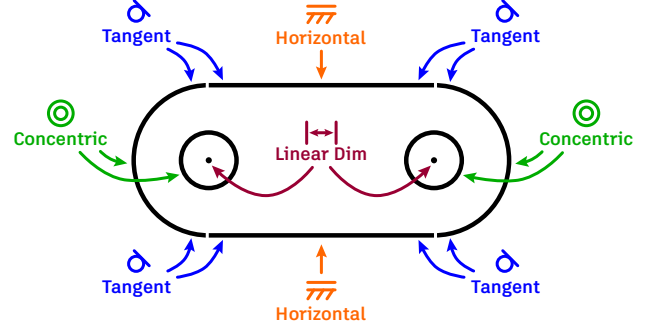


Figure A.1. An example sketch illustrating how constraints and dimensions reference geometric primitives such as points, lines, arcs, and circles. A constraint solver enforces these relationships, ensuring that a change in one parameter propagates consistently throughout the sketch.

A.3. Sketch Tokenization

Our tokenization of sketches defines a diverse vocabulary of token types to represent the heterogeneous elements of a sketch. There are distinct token categories for primitive types, constraint and dimension types, and special markers (e.g., `<SOS>`, `<EOS>`, `<PAD>`). In our approach, constraint tokens, dimension tokens, and primitive reference tokens are the primary outputs of the model. These tokens are strictly categorical, reflecting the discrete nature of constraint types and their relationship to previously defined primitives. For example, a perpendicular constraint might be tokenized as `(<PER>, <REF_A>, <REF_B>)`, where `<REF_A>` and `<REF_B>` are reference tokens pointing to two lines introduced earlier in the sequence.

While geometric primitives also contain continuous parameters (coordinates, radii, angles, etc.), these parameters are not predicted by our model. Instead, they are treated as input to inform constraint generation. To incorporate this information, each primitive’s continuous parameters are embedded in a separate stream of tokens for input only. The generative process focuses on discrete constraints and dimensions that reference the primitives, leaving numeric values for dimensions to be resolved by the constraint solver. This design choice leverages the solver’s robust capacity to converge on valid parameter assignments, allowing the model to prioritize structural correctness and alignment with design objectives.

A.4. SketchGraphs Dataset

While the main body describes how the SketchGraphs dataset was filtered and converted, additional details regarding the motivation and practical considerations of each step are provided here. The primary goal of these refinements is to produce a clean, representative subset of sketches and

ensure each example aligns with standard engineering constraints.

In Table A.1 we list out the supported constraint and dimension types in Onshape terminology that we included in the training data. Notably, we filter out less prevalent constraints (Symmetric, Normal, Pattern) and dimensions (CenterLine, Projected) to focus the learning task on the core geometric relationship types which form the backbone of sketch geometry. By removing these non-core constraints, we simplify the constraint vocabulary the model must learn while still covering the vast majority of design intent in sketches.

Table A.1. Supported Constraints and Dimensions

Constraints	Dimensions
Coincident	Diameter
Horizontal	Radius
Vertical	Distance
Parallel	Angle
Perpendicular	Length
Tangent	
Midpoint	
Equal	
Offset	
Concentric	

We next eliminate redundant constraints by deduplicating overlapping coincident points. We identify groups of points that all coincide and merge or remove duplicate coincident constraints among them. This deduplication of coincident points removes unnecessary edges in the constraint graph, reducing its complexity without altering the sketch’s geometry. This focuses the model on the unique geometric relationships and avoids penalizing it for not outputting repetitive constraints that do not add new information. To avoid bias from repeated structures, we also deduplicate very similar or identical sketches in the dataset. We detect and remove duplicate sketches so that each unique sketch structure is represented more evenly.

After applying the above filters, we verify each sketch’s constraints for solver solvability. Any sketch that the solver identifies as unsolvable is removed from the training set for the SFT model training. This step guarantees that the model trains only on valid, feasible sketches that correspond to a realizable geometry. We also exclude sketches that are grossly under-constrained, where the solver indicates many degrees of freedom remain, since they may not demonstrate clear constraint interactions for the model to learn. However, we add these sketches back for model fine-tuning.

Finally, we fix at least one point in each sketch to lock its position. Because the SketchGraphs data often provides no absolute anchor in the plane, many sketches exhibit degrees of freedom that allow global translation or rotation without

altering constraints internally. In a typical design environment, at least one point or an entire component is fixed to serve as a reference. Fixing a point eliminates global translational and rotational degrees of freedom, effectively locking the sketch in a consistent pose.

Table A.2 provides detailed statistics of the SketchGraphs dataset after preprocessing. At the dataset level, the resulting set contains approximately 2.8 million sketches. Among these, only 8.27% of sketches are fully-constrained (FC), highlighting the rarity of sketches that require no additional constraints. Around 16.11% are over-constrained (OC), while 1.62% are unsolvable. A majority (93.70%) of sketches are stable when stability is evaluated using a 4-bin discretization of geometry positions.

At the sketch level, the average sketch consists of about 15 geometric entities and contains roughly 7 constraints and 1 dimension, although there is considerable variation (standard deviation 7.27, 5.48, and 1.67, respectively). Additionally, point-level and curve-level fully-constrained percentages per sketch average at approximately 27% and 33%, respectively, indicating that most sketches are significantly under-constrained at the primitive level.

B. Architecture and Experiment Details

We discuss additional details regarding the model architecture, training, and experiments.

B.1. Vitruvion

We use Vitruvion as the core constraint generation model for all post-training algorithms. Our implementation is adapted to work with the Fusion sketch representation, which treats all points as distinct geometric primitives. This differs from Onshape, which introduces the concept of sub-primitives – geometric entities can own points (e.g., a line owns its start and end points). In the tokenized geome-

Table A.2. Statistics of the SketchGraphs dataset after preprocessing.

Dataset-Level Statistics				
Sketch Count	2,784,964			
% FC	8.27	% Not Solvable	1.62	
% OC	16.11	% Stable (bins=4)	93.70	
Sketch-Level Statistics				
	Mean \pm Std	Min	Median	Max
Entity Count	14.68 \pm 7.27	1	13	64
Constraint Count	6.53 \pm 5.48	0	5	52
Dimension Count	1.08 \pm 1.67	0	0	42
% Point FC	27.13 \pm 22.72	0.00	20.00	100.00
% Curve FC	33.48 \pm 29.49	0.00	28.57	100.00

try sequence, each geometric entity is represented by its top-level primitive along with a nested list of its associated sub-primitives. The pointer network can then reference both sub-primitives and standard primitives within the index space of the tokenized geometry sequence. By contrast, in Fusion there is no concept of “sub-primitives” – all indices in the tokenized geometry sequence are associated with independent primitives. When pre-processing the data, we combine duplicate points in the SketchGraphs data and initialize these as separate points (i.e. not owned by a curve).

We additionally include a learned embedding for each entity indicating whether or not the entity is fixed or not. As mentioned in Appendix A.4, at least one fixed entity is necessary to act as an anchor to the rest of the sketch. In order for an entity to be fully constrained, the constraint graph must connect to a fixed entity. We posit that this information is valuable for the task of fully constraining sketches.

Our implementation represents curves, circles, and arcs using 5 points extracted along the path of the shape. This differs from Vitruvion which uses the parameters of the shape such as start/end points, center, radius, and arc midpoint. Lastly, we model constraints using the given (user) order rather than ordering based on the referenced primitives.

B.2. Preference-based Optimization Algorithms

The hyperparameters for our preference-based optimization algorithms are presented in Table B.1. Both DPO and Expert Iteration (ExIt) methods are initialized from the SFT model and undergo 2 full rounds of data generation using a temperature of 1.0 followed by policy improvement. The DPO implementation has additional hyper-parameters: a β parameter controls preference strength, a small SFT loss weight combines the DPO loss with a standard cross-entropy loss on the positive sample τ_w , and a label smoothing weight reduces model overconfidence. These settings were determined through preliminary experiments to optimize model performance.

Table B.1. Training hyperparameters for preference-based optimization algorithms.

Hyperparameters	ExIt	DPO
Batch size	64	64
Rounds (N)	2	2
Learning rate	1e-6	1e-5
Sampling temperature (data)	1.0	1.0
β (DPO)	-	0.1
SFT weight	-	0.05
Label smoothing weight	-	0.3

In the data generation phase, ExIt uses rejection sampling to filter out any under-constrained, over-constrained,

or unsolvable model outputs. For DPO, we find all pairs (τ_w, τ_l) of model outputs for the same sketch where τ_w is fully-constrained and τ_l is under-constrained, over-constrained, or unsolvable. In order to help DPO better distinguish between the positive and negative examples, we limit τ_l to have less than 90% fully constrained curves.

B.3. RL algorithms

For the rewards, we used $r_{\text{unstable}} = -0.25$ as a penalty for unstable sketches, $r_{\text{NS}} = -1.0$ as a penalty for not solvable sketches, $r_{\text{OC}} = -1.0$ as a penalty for over-constrained sketches, and $r_{\text{F}} = -0.5$ as a penalty for sketches resulting in other failures. Other training hyperparameter choices are shown in Table B.2.

Table B.2. Training hyperparameters for RL algorithms.

Hyperparameters	ReMax	RLOO	GRPO
Batch size	32	32	32
Group sample size	-	8	8
Learning rate	1e-5	1e-5	1e-5
Sampling temperature	1.0	1.0	1.0
Reference update timesteps	100	100	100
KL penalty added to rewards	0.01	0.01	0.0
KL regularization β	-	-	0.01
Policy clipping threshold ϵ	-	-	0.2

C. Additional results

C.1. Diversity

Table C.1 presents the diversity metrics for constraint generation across different models. The Vitruvion base model demonstrates the highest diversity with 65.23% unique generations and a relatively low Mean Intersection over Union (MIoU) of 0.623, indicating substantial variation between generated constraints. In contrast, RLOO and GRPO show the least diversity, with 32.11% and 33.95% unique sketches respectively, and high MIoU values exceeding 0.88, suggesting considerable overlap in their generations. Expert Iteration achieves a better balance, maintaining relatively high diversity (62.80% unique) while improving on the base model’s performance. Standard SFT and Iterative DPO fall between these extremes, with the latter showing moderately improved diversity metrics over SFT.

C.2. Number of DPO/ExIt Iterations

Figure C.1 shows the performance of the preference-based optimization algorithms across training rounds. Expert iteration shows better performance at generating fully-constrained and not over-constraining sketches compared to DPO. One possible reason for this is that the process of selecting positive/negative example pairs for DPO is

Table C.1. Diversity results computed across 8 generations per sketch. Unique@8 is the percentage of the time that the model generates a unique set of constraints for each sketch, compared to the other generations for the same sketch. We measure uniqueness with the Weisfeiler Lehman (WL) graph hash with 4 quantization bins. MIoU is the average intersection over union of the generated constraints between the other generations for each sketch.

Model	% Unique@8 \uparrow	MIoU@8 \downarrow
Vitruvion (base)	65.23	0.623
SFT	46.71	0.782
Iterative DPO	52.79	0.775
Expert Iteration	62.80	0.720
ReMax	35.80	0.877
RLOO	32.11	0.892
GRPO	33.95	0.881

more restrictive since each positive (fully-constrained) example must be matched with an under-constrained or over-constrained example for the same sketch.

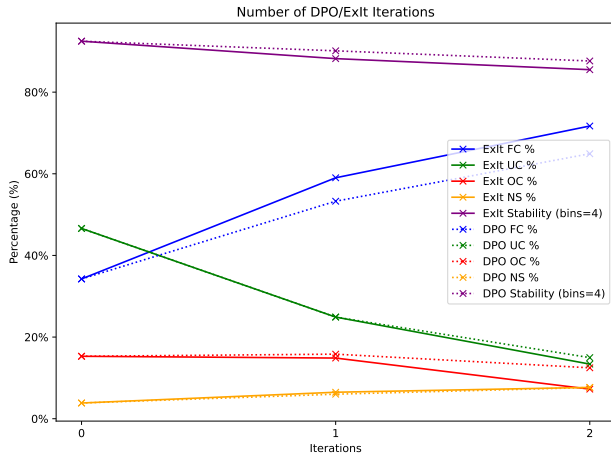


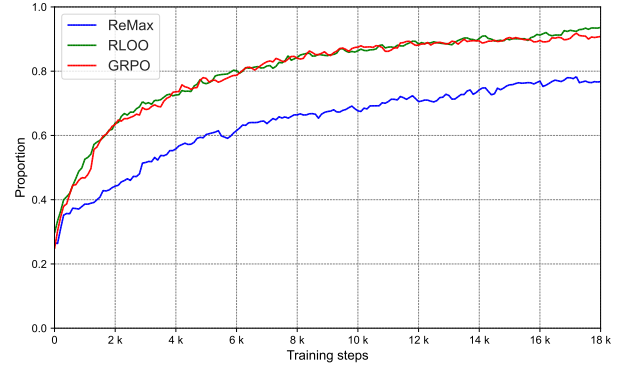
Figure C.1. Performance across rounds for Iterative DPO and Expert Iteration. Results are the mean of $K = 8$ samples. The initial model at $t = 0$ is the SFT model

C.3. RL Training curves

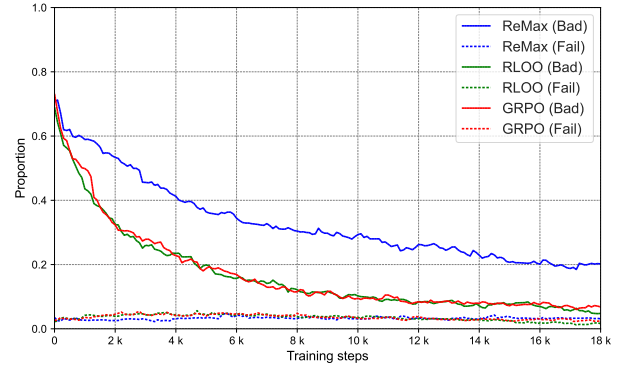
Figure C.2 shows training performance over time for the online reinforcement learning algorithms.

C.4. Failed Attempts

Despite our efforts to leverage reinforcement learning for constraint generation, we encountered several dead ends. Each failed attempt underlines a fundamental challenge in aligning reward signals and exploration strategies with the requirements of geometric constraint generation. Below, we discuss three key failures, followed by brief summaries of the lessons drawn from each.



(a) Proportion of successfully constrained sketches



(b) Proportion of unsuccessfully constrained sketches

Figure C.2. Proportion of (a) successful sketches (Fully-Constrained and not Over-Constrained), and (b) badly-constrained sketches (either Under-Constrained and/or Over-Constrained) and failed sketches (resulting in constraints solver error) over training for RL methods (ReMax, RLOO and GRPO). (Stability is not considered here in determining successfulness.)

C.4.1. PPO with a Learned Reward Model

We first attempted to train a policy using PPO, guided by a learned reward model predicting how well the generated constraints would align with desired outcomes. This reward model serves as a surrogate model of the constraint solver, estimating the curve and point fully-constrained percentage, fully-constrained and under-constrained status, and stability. Unfortunately, the agent over-fit the reward model’s idiosyncrasies instead of genuinely improving constraint quality. In our case, PPO steadily increased the reward model’s score, but the rate of curve fully-constrained percentage actually dropped, which is evident that the policy was “reward hacking” the learned metric.

Several practical issues led to this failure. First, we lack diverse training samples for the reward model, especially for over-constrained or edge-case scenarios. The reward model was trained on two different settings, either on sparse per-sequence labels (only knowing the true evaluation met-

rics given an entire constraint set) or on per-constraint feedback. Both schemes suffered from limited coverage of failure modes. When PPO began producing novel constraint combinations outside the training distribution, the reward model was out of its depth. In our implementation, the reward model remained fixed during PPO fine-tuning; as the policy explored new regions of the constraint space, the frozen reward model’s prediction errors grew unchecked.

C.4.2. PPO with Solver-based Rewards

Another approach replaced the learned reward model with direct solver feedback, providing a reward only when the entire constraint sequence is generated. Although this feedback was unambiguously correct, it proved extremely sparse, the distribution of rewards remained highly skewed, with most episodes clustered near the lower or neutral end and only infrequent high-reward successes, causing training to collapse. For the policy gradient approach, such sporadic positive returns can still nudge the policy upward in proportion to the log probability of successful episodes. In contrast, the PPO algorithm sees little incremental feedback to guide learning, sudden high rewards are either clipped or overshadowed by large variance in advantage estimates.

C.4.3. Logic-based Action Masking

Finally, we tested logit masking to disallow certain “invalid” actions. In principle, this was meant to help by preventing the agent from exploring blatantly wrong moves. Surprisingly, this logit masking made learning worse for all our RL algorithms. One theoretical reason is that the mask, while eliminating invalid actions, also over-constrained the policy’s exploration. Contrary to expectations, blocking these actions harmed training. By never letting the agent attempt blatantly invalid moves, the model lost valuable negative feedback signals and drastically curtailed exploration. Another theoretical concern is that dynamic action masking can complicate the Markov Decision Process. So the issue is likely not that the concept of masking is invalid, but rather that it altered the learning dynamics in our specific setting.