# Imagine This! Scripts to Compositions to Videos

Tanmay Gupta[1], Dustin Schwenk[2], Ali Farhadi[2,3],
Derek Hoiem[1], and Aniruddha Kembhavi[2]

[1] University of Illinois Urbana-Champaign
[2] Allen Institute for Artificial Intelligence
[3] University of Washington

**Abstract.** Imagining a scene described in natural language with realistic layout and appearance of entities is the ultimate test of spatial, visual, and semantic world knowledge. Towards this goal, we present the **C**omposition, **R**etrieval **a**nd **F**usion Ne**t**work (CRAFT), a model capable of learning this knowledge from video-caption data and applying it while generating videos from novel captions. CRAFT explicitly predicts a temporal-layout of mentioned entities (characters and objects), retrieves spatio-temporal entity segments from a video database and fuses them to generate scene videos. Our contributions include sequential training of components of CRAFT while *jointly* modeling layout and appearances, and losses that encourage learning compositional representations for retrieval. We evaluate CRAFT on *semantic fidelity* to caption, *composition consistency*, and *visual quality*. CRAFT outperforms direct pixel generation approaches and generalizes well to unseen captions and to unseen video databases with no text annotations. We demonstrate CRAFT on FLINTSTONES, a new richly annotated video-caption dataset with over 25000 videos. For a glimpse of videos generated by CRAFT, see https://youtu.be/688Vv86n0z8.
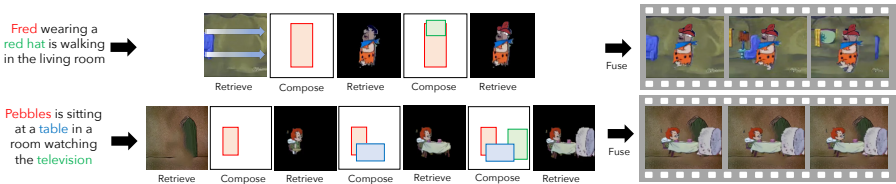
**Fig. 1.** Given a novel description, CRAFT sequentially composes a scene layout and retrieves entities from a video database to create complex scene videos.

## 1 Introduction

Consider the scene description: *Fred is wearing a blue hat and talking to Wilma in the living room. Wilma then sits down on a couch.* Picturing the scene in our mind requires the knowledge of plausible locations, appearances, actions, and interactions of characters and objects being described, as well as an ability to understand and translate the natural language description into a plausible visual instantiation. In this work, we introduce Semantic Scene Generation (SSG), the task of generating complex scene videos from rich natural language descriptions which requires jointly modeling the layout and appearances of entities mentioned in the description. SSG models are trained using a densely annotated video

dataset with scene descriptions and entity bounding boxes. During inference, the models must generate videos for novel descriptions (unseen during training).

Modelling the layout and appearances of entities for descriptions like the one above poses several challenges: (a) **Entity Recall** - the video must contain the relevant characters (Fred, Wilma), objects (blue hat, couch) and background (setting that resembles a living room); (b) **Layout Feasibility** - characters and objects must be placed at plausible locations and scales (Fred, Wilma and the couch should be placed on the ground plane, the hat must lie on top of Fred's head); (c) **Appearance Fidelity** - entity appearance, which may be affected by identity, pose, action, attributes and layout, should respect the scene description; (d) **Interaction Consistency** - appearance of characters and objects must be *consistent with each other* given the described, sometimes implicit, interaction (Fred and Wilma should face each other as do people when they talk to each other); (f) **Language Understanding** - the system must be able to understand and translate a natural language description into a plausible visual instantiation.

Currently, the dominant approaches to conditional generation of visual data from text rely on directly learning distributions in a *high dimensional pixel space*. While these approaches have shown impressive results for aligned images of objects (faces, birds, flowers, etc.), they are often inadequate for addressing the above challenges, due to the *combinatorial explosion* of the image space arising from multiple characters and objects with significant appearance variations arranged in a large number of possible layouts. In contrast, our proposed **C**omposition, **R**etrieval **a**nd **F**usion Ne**t**work (CRAFT) explicitly models the spatio-temporal layout of characters and objects in the scene jointly with entity appearances. Unlike pixel generation approaches, our appearance model is based on text to entity segment retrieval from a video database. Spatio-temporal segments are extracted from the retrieved videos and fused together to generate the final video. The layout composition and entity retrieval work in a sequential manner which is determined by the language input. Factorization of our model into composition and retrieval stages alleviates the need to directly model pixel spaces, results in an architecture that more easily exploits location and appearance contextual cues, and renders a more interpretable output.

Towards the goal of SSG, we introduce FLINTSTONES, a densely annotated dataset based on *The Flintstones* animated series, consisting of over 25000 videos, each 75 frames long. FLINTSTONES has several advantages over using a random sample of internet videos. First, in a closed world setting such as a television series, the most frequent characters are present in a wide variety of settings, which serves as a more manageable learning problem than a sparse set obtained in an open world setting. Second, the flat textures in animations are easier to model than real world videos. Third, in comparison to other animated series, The Flintstones has a good balance between having fairly complex interactions between characters and objects while not having overly complicated, cluttered scenes. For these reasons, we believe that the FLINTSTONES dataset is semantically rich, preserves all the challenges of text to scene generation and is a good stepping stone towards real videos. FLINTSTONES consists of an 80-10-10

train-val-test split. The train and val sets are used for learning and model selection respectively. Test captions serve as novel descriptions to generate videos at test time. To quantitatively evaluate our model, we use two sets of metrics. The first measures *semantic fidelity* of the generated video to the desired description using entity noun, adjective, and verb recalls. The second measures *composition consistency*, *i.e.* the consistency of the appearances, poses and layouts of entities with respect to other entities in the video and the background.

We use FLINTSTONES to evaluate CRAFT and provide a detailed ablation analysis. CRAFT outperforms baselines that generate pixels directly from captions as well as a whole video retrieval approach (as opposed to modeling entities). It generalizes well to unseen captions as well as unseen videos in the target database. Our quantitative and qualitative results show that for simpler descriptions, CRAFT exploits location and appearance contextual cues and outputs videos that have consistent layouts and appearances of described entities. However, there is tremendous scope for improvement. CRAFT can fail catastrophically for complex descriptions (containing large number of entities, specially infrequent ones). The adjective and verb recalls are also fairly low. We believe SSG on FLINTSTONES presents a challenging problem for future research.

## 2  Related Work

**Generative models** Following pioneering work on Variational Autoencoders [1] and Generative Adversarial Networks [2], there has been tremendous interest in generative modelling of visual data in a high dimensional pixel space. Early approaches focused on unconditional generation [3–6], whereas recent works have explored models conditioned on simple textual inputs describing objects [7–11]. While the visual quality of images generated by these models has been steadily improving [12, 13], success stories have been limited to generating images of aligned objects (e.g. faces, birds, flowers), often training one model per object class. In contrast, our work deals with generating complex scenes which requires modelling the layout and appearances of multiple entities in the scene.

Of particular relevance is the work by Hong *et al.* [14] who first generate a coarse semantic layout of bounding boxes, refine that to segmentation masks and then generate an image using an image-to-image translation model [15,16]. A limitation of this approach is that it assumes a fixed number of object classes (80 in their experiments) and struggles with the usual challenge of modeling high dimensional pixel spaces such as generating coherent entities. Formulating appearance generation in terms of entity retrieval from a database allows our model to scale to a large number of entity categories, guarantee intra-entity coherence and allows us to focus on the semantic aspects of scene generation and inter-entity consistency. The retrieval approach also lends itself to generating videos without significant modification. There have been attempts at extending GANs for unconditional [17,18] as well as text conditional [19] video generation, but quality of generated videos is usually worse than that of GAN generated images unless used in very restrictive settings. A relevant generative modelling approach is by Kwak *et al.* [20] who proposed a model in which parts of the
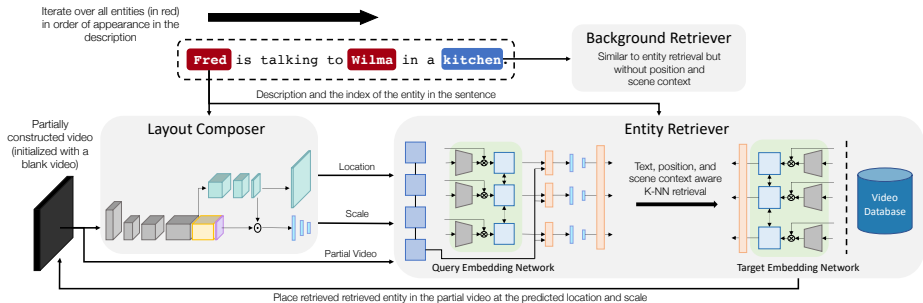
**Fig. 2. Overview** of **C**omposition, **R**etrieval **a**nd **F**usion Ne**t**work (CRAFT), consisting of three parts: *Layout Composer*, *Entity Retriever* and *Background Retriever*. CRAFT begins with an empty video and sequentially adds entities mentioned in the input description at locations and scales predicted by the Layout Composer.

image are generated sequentially and combined using alpha blending. However, this work does not condition on text and has not been demonstrated on complex scenes. Another relevant body of work is by Zitnick *et al.* [21–23] who compose static images from descriptions with clipart images using a Conditional Random Field formulation.

To control the structure of the output image, a growing body of literature conditions image generation on a wide variety of inputs ranging from keypoints [24] and sketches [25] to semantic segmentation maps [15]. In contrast to these approaches which condition on *provided* location, our model *generates* a plausible scene layout and then conditions entity retrieval on this layout.

**Phrase Grounding and Caption-Image Retrieval.** The entity retriever in CRAFT is related to caption based image retrieval models. The caption-image embedding space is typically learned by minimizing a ranking loss such as a triplet loss [26, 26–29]. Phrase grounding [30] is another closely related task where the goal is to localize a region in an image described by a phrase.

One of our contributions is enriching the semantics of embeddings learned through triplet loss by simultaneously minimizing an auxiliary classification loss based on noun, adjective and verb words associated with an entity in the text description. This is similar in principle to [31] where auxiliary autoencoding losses were used in addition to a primary binary prediction loss to learn robust visual semantic embeddings. Learning shared representations across multiple related tasks is a key concept in multitask learning [32, 33].

## 3  Model

Figure 2 presents an overview of **C**omposition, **R**etrieval **a**nd **F**usion Ne**t**work which consists of three parts: *Layout Composer*, *Entity Retriever*, and *Background Retriever*. Each is a neural network that is trained independently using ground truth supervision. At inference time, CRAFT begins with an empty video and adds entities in the scene sequentially based on the order of appearance in the description. At each step, the *Layout Composer* predicts a location and scale for an entity given the text and the video constructed so far. Then, conditioned

on the predicted location, text, and the partially constructed video, the *Entity Retriever* produces a query embedding that is looked up against the embeddings of entities in the target video database. The entity is cropped from the retrieved video and placed at the predicted location and scale in the video being generated. Alternating between the *Layout Composer* and *Entity Retriever* allows the model to condition the layout of entities on the appearance and vice versa. Similar to *Entity Retriever*, the *Background Retriever* produces a query embedding for the desired scene from text and retrieves the closest background video from the target database. The retrieved spatio-temporal entity segments and background are fused to generate the final video. We now present the notation used in the rest of the paper, followed by architecture and training details for the three components.

| Caption | |
|---|---|
| $T$ | Caption with length $|T|$ |
| $\{E_i\}_{i=1}^n$ | $n$ entities in $T$ in order of appearance |
| $\{e_i\}_{i=1}^n$ | entity noun positions in $T$ |
| **Video** | |
| $F$ | number of frames in a video |
| $\{(l_i, s_i)\}_{i=1}^n$ | position of entities in the video |
| $l_i$ | entity bounding box at each frame ($\{(x_{if}, y_{if}, w_{if}, h_{if})\}_{f=1}^F$) |
| $s_i$ | entity pixel segmentation mask at each frame |
| $V_{i-1}$ | partially constructed video with entities $\{E_j\}_{j=1}^{i-1}$ |
| $V (= V_n)$ | full video containing all entities |
| $\{(V^{[m]}, T^{[m]})\}_{m=1}^M$ | training data points, where $M$ = number of data points |

## 3.1 Layout Composer

The layout composer is responsible for generating a plausible layout of the scene consisting of the locations and scales of each character and object mentioned in the scene description. Jointly modeling the locations of all entities in a scene presents fundamentally unique challenges for spatial knowledge representation beyond existing language-guided localization tasks. Predicting plausible locations and scales for objects not yet in an image requires a significant amount of *spatial knowledge* about people and objects, in contrast to text based object localization which relies heavily on appearance cues. This includes knowledge like – a hat goes on top of a person's head, a couch goes under the person sitting on it, a person being talked to faces the person speaking instead of facing away, tables are short and wide while standing people are tall and thin, etc.

Figure 3 presents a schematic for the layout composer. Given the varying number of entities across videos, the layout composer is setup to run in a sequential manner over the set of distinct entities mentioned in a given description. At each step, a text embedding of the desired entity along with a partially constructed video (consisting of entities fused into the video at previous steps) are input to the model which predicts distributions for the location and scale of the desired entity.

The layout composer models $P(l_i|V_{i-1}, T, e_i; \theta_{loc}, \theta_{sc})$, the conditional distribution of the location and scale (width and height normalized by image size)
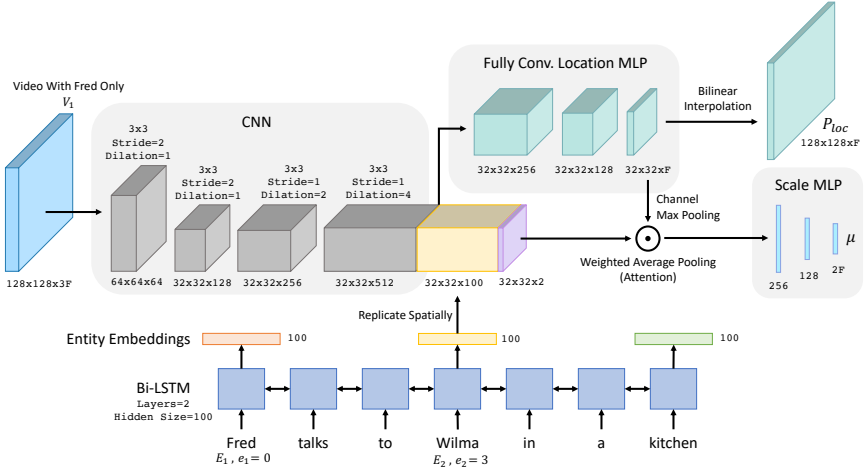
**Fig. 3. Layout Composer** is run sequentially through the set of entities in the description, predicting the distributions for the location and scale of the desired entity.

of the $i^{th}$ entity given the text, entity noun position in tokenized text, and the partial video with previous entities. Let $C_i$ denote the conditioning information, $(V_{i-1}, T, e_i)$. We factorize the position distribution into location and scale components as follows:

$$P(l_i|C_i; \theta_{loc}, \theta_{sc}) = \prod_{f=1}^{F} P_{loc}^f(x_{if}, y_{if}|C_i; \theta_{loc}^f) \cdot P_{sc}^f(w_{if}, h_{if}|x_{if}, y_{if}, C_i; \theta_{sc}^f) \quad (1)$$

$\theta_{loc} = \{\theta_{loc}^f\}_{f=1}^F$ and $\theta_{sc} = \{\theta_{sc}^f\}_{f=1}^F$ are learnable parameters. $P_{loc}^f$ is modelled using a network that takes $C_i$ as input and produces a distribution over all pixel locations for the $f^{th}$ image frame. We model $P_{sc}^f$ using a Gaussian distribution whose mean $\mu_f$ and covariance $\Sigma_f$ are predicted by a network given $(x_i, y_i, C_i)$. Parameters $\theta_{loc}$ and $\theta_{sc}$ are learned from ground truth position annotations by minimizing the following maximum likelihood estimation loss:

$$\sum_{m=1}^{M} \sum_{i=1}^{n^{[m]}} \sum_{f=1}^{F} \Big[ -\log(P_{loc}^f(x_{if}^{[m]}, y_{if}^{[m]}|C_i^{[m]}; \theta_{loc}^f)) + 0.5 \cdot \log(\det(\Sigma(x_{if}, y_{if}, C_i; \theta_{sc}^f))) +$$

$$0.5 \cdot (z_{if}^{[m]} - \mu_f(D_i^{[m]}; \theta_{sc}^f))^T \Sigma_f^{-1}(z_{if}^{[m]} - \mu_f(D_i^{[m]}; \theta_{sc}^f)) + \log(2\pi) \Big] \quad (2)$$

where $z_{if} = [w_{if}; h_{if}]$ & $D_i^{[m]} = (x_i^{[m]}, y_i^{[m]}, C_i^{[m]})$. For simplicity, we manually set and freeze $\Sigma$ to an isometric diagonal covariance matrix with variance of 0.005.

**Feature Computation Backbone.** The location and scale predictors have an identical feature computation backbone comprising of a CNN and a bidirectional LSTM. The CNN encodes $V_{i-1}$ (8 sub-sampled frames concatenated along the channel dimension) as a set of convolutional feature maps which capture appearance and positions of previous entities in the scene. The LSTM is used to encode the entity $E_i$ for which the prediction is to be made along with semantic context available in the caption. The caption is fed into the LSTM and the hidden output at $e_i^{th}$ word position is extracted as the entity text encoding. The text encoding is replicated spatially and concatenated with convolutional features and 2-D

grid coordinates to create a representation for each location in the convolutional feature grid that is aware of visual, spatial, temporal, and semantic context.

**Location Predictor.** $P_{loc}^f$ is modelled using a Multi Layer Perceptron (MLP) that produces a score for each location. This map is bilinearly upsampled to the size of input video frames. Then, a softmax layer over all pixels produces $P_{loc}^f(x, y|C; \theta_{loc}^f)$ for every pixel location $(x, y)$ in the $f^{th}$ video frame.

**Scale Predictor.** Features computed by the backbone at a particular $(x, y)$ location are selected and fed into the scale MLP that produces $\mu_f(x_i, y_i, C_i; \theta_{sc}^f)$.

**Feature sharing and multitask training.** While it is possible to train a separate network for each $\{P_{loc}^f, \mu_f\}_{f=1}^F$, we present a pragmatic way of sharing features and computation for different frames and also between the location and scale networks. To share features and computation across frames, the location network produces $F$ probability maps in a single forward pass. This is equivalent to sharing all layers across all $P_{loc}^f$ nets except for the last layer of the MLP that produces location scores. Similarly, all the $\mu_f$ nets are also combined into a single network. We refer to the combined networks by $P_{loc}$ and $\mu$.

In addition, we also share features across the location and scale networks. First, we share the feature computation backbone, the output from which is then passed into location and scale specific layers. Second, we use a soft-attention mechanism to select likely positions for feeding into the scale layers. This conditions the scale prediction on the plausible locations of the entity. We combine the $F$ spatial maps into a single attention map through max pooling. This attention map is used to perform weighted average pooling on backbone features and then fed into the scale MLP. Note that this is a differentiable greedy approximation to find the most likely location (by taking argmax of spatial probability maps) and scale (directly using output of $\mu$, the mode for a gaussian distribution) in a single forward pass. To keep training consistent with inference, we use the soft-attention mechanism instead of feeding ground-truth locations into $\mu$.

### 3.2   Entity Retriever

The task of the entity retriever is to find a spatio-temporal patch within a target database that matches an entity in the description *and* is consistent with the video constructed thus far – the video with all previous entities retrieved and placed in the locations predicted by the layout network. We adopt an embedding based lookup approach for entity retrieval. This presents several challenges beyond traditional image retrieval tasks. Not only does the retrieved entity need to match the semantics of the description but it also needs to respect the implicit relational constraints or context imposed by the appearance and locations of other entities. E.g. for *Fred is talking to Wilma*, it is not sufficient to retrieve *a Wilma*, but one who is also facing in the right direction, i.e. towards *Fred*.

The Entity Retriever is shown in Figure 4 and consists of two parts: (i) query embedding network $Q$, and (ii) target embedding network $R$. $Q$ and $R$ are learned using the query-target pairs $\langle (T^{[m]}, e_i^{[m]}, l_i^{[m]}, V_{i-1}^{[m]}), (V^{[m]}, l_i^{[m]}, s_i^{[m]}) \rangle_{i,m}$ in
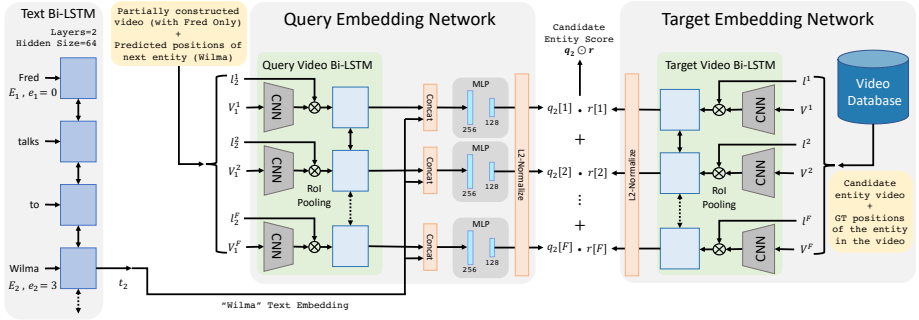
**Fig. 4. Entity Retriever** retrieves spatio-temporal patches from a target database that match entity description as encoded by the query embedding network.

the training data. For clarity, we abbreviate $Q(T^{[m]}, e_i^{[m]}, l_i^{[m]}, V_{i-1}^{[m]})$ as $q_i^{[m]}$ and $R(V^{[m]}, l_i^{[m]}, s_i^{[m]})$ as $r_i^{[m]}$. At each training iteration, we sample a mini-batch of $B$ pairs without replacement and compute embeddings $\{(q_{i_b}^{[m_b]}, r_{i_b}^{[m_b]})\}_{b=1}^{B}$ where $q$ and $r$ are each sequence of $F$ embeddings corresponding to $F$ video frames. The model is trained using a triplet loss computed on *all possible* triplets in the mini-batch. Let $\delta_b$ denote the set of all indices from 1 to $B$ except $b$. The loss can then be defined as

$$\mathcal{L}_{triplet} = \frac{1}{B \cdot (B-1)} \sum_{b=1}^{B} \sum_{b^- \in \delta_b} \Big[ \max(0, \gamma + q_{i_b}^{[m_b]} \odot r_{i_{b^-}}^{[m_{b^-}]} - q_{i_b}^{[m_b]} \odot r_{i_b}^{[m_b]}) +$$

$$\max(0, \gamma + q_{i_{b^-}}^{[m_{b^-}]} \odot r_{i_b}^{[m_b]} - q_{i_b}^{[m_b]} \odot r_{i_b}^{[m_b]}) \Big] \quad (3)$$

where $q \odot r = \frac{1}{F} \sum_{f=1}^{F} q[f] \cdot r[f]$ is the average dot product between corresponding query and target frame embeddings. We use a margin of $\gamma = 0.1$.

**Auxiliary Multi-label Classification Loss** We found that models trained using triplet loss alone could simply learn a one-to-one mapping between ground truth text and entity video segments with poor generalization to unseen captions and database videos. To guide the learning to utilize the compositional nature of text and improve generalization, we add an auxiliary classification loss on top of the embeddings. The key idea is to enrich the semantics of the embedding vectors by predicting the noun, adjectives, and action words directly associated with the entity in the description. For example, *Wilma*'s embedding produced by the query and target embedding networks in *Fred is talking to a happy Wilma who is sitting on a chair.* is forced to predict *Wilma*, *happy* and *sitting* ensuring their representation in the embeddings. A vocabulary $\mathcal{W}$ is constructed of all nouns, adjectives and verbs appearing in the training data. Then for each sample in the mini-batch, an MLP is used as a multi-label classifier to predict associated words from the query and target embeddings. Note that a single MLP is used to make these noun, adjective and verb predictions on *both* query and target embeddings.

**Query Embedding Network ($Q$).** Similar to the layout composer's feature computation backbone, $Q$ consists of a CNN to independently encode every frame of $V_{i-1}$ and an LSTM to encode $(T, e_i)$ which are concatenated together

along with a 2-D coordinate grid to get per-frame feature maps. However, unlike layout composer, the query embedding network also needs to be conditioned on the position $l_i$ where entity $E_i$ is to be inserted in $V_{i-1}$. To get location and scale specific query embeddings, we use a simplified RoIAlign (RoIPool with RoI quantization and bilinear interpolation) mechanism to crop out the per-frame feature maps using the corresponding bounding box $l_i^f$ and scaling it to a $7 \times 7$ receptive field. The RoIAlign features are then averaged along the spatial dimensions to get the vector representations for each time step independently. An LSTM applied over the sequence of these embeddings is used to capture temporal context. The hidden output of the LSTM at each time step is normalized and used as the frame query embedding $q[f]$.

**Target Embedding Network ($R$).** Since during inference, $R$ needs to embed entities in the target database which do not have text annotations, it does not use $T$ as an input. Thus, $R$ is similar to $Q$ but without the LSTM to encode the text. In our experiments we found that using 2-D coordinate features in both query and target networks made the network susceptible to ignoring all other features since it provides an easy signal for matching ground truth query-target pairs during training. This in turn leads to poor generalization. Thus, $R$ has no 2-D coordinate features.

### 3.3   Background Retriever

The task of the background retriever is to find a background scene that matches the setting described in the description. To construct a database of backgrounds without characters in them, we remove characters from videos (given bounding boxes) and perform hole filling using PatchMatch [34]. The background retriever model is similar to the entity retriever with two main differences. First, since the whole background scene is retrieved instead of entity segments, the conditioning on position is removed from both query and database embedding networks replacing RoI pooling with global average pooling. Second, while ideally we would like scene and entity retrieval to be conditioned on each other, for simplicity we leave this to future work and currently treat them independently. These modifications essentially reduce the query embedding network to a text Bi-LSTM whose output at the background word location in the description is used as the query embedding, and the target embedding network to a video Bi-LSTM without RoI pooling. The model is trained using just the triplet loss.

## 4   The Flintstones Dataset

**Composition.** The FLINTSTONES dataset is composed of 25184 densely annotated video clips derived from the animated sitcom *The Flintstones*. Clips are chosen to be 3 seconds (75 frames) long to capture relatively small action sequences, limit the number of sentences needed to describe them and avoid scene and shot changes. As shown in Figure 5, annotations contain clip's characters,
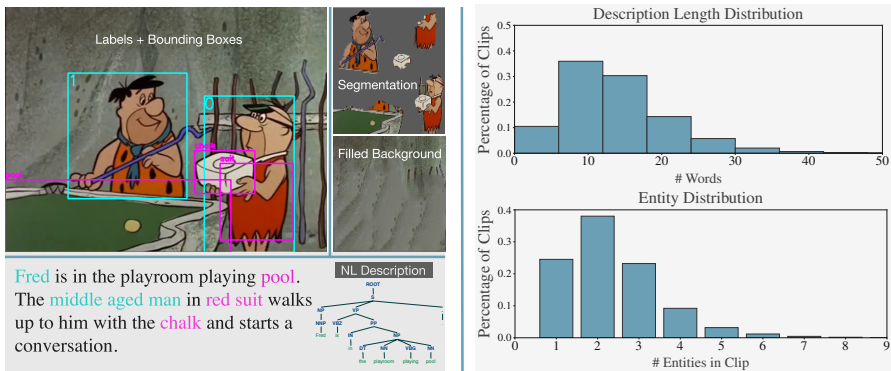
**Fig. 5.** FLINTSTONES **dataset.** Left: Overview of clip annotations. Right: Distribution of description lengths and entity counts in the dataset.

setting, and objects being interacted with marked in text as well as their bounding boxes in all frames. FLINTSTONES has a 80-10-10 train-val-test split.

**Clip Annotation.** Dense annotations are obtained in a multi-step process: identification and localization of characters in keyframes, identification of the scene setting, scene captioning, object annotation, and entity tracking to provide annotations for all frames. The dataset also contains segmentation masks for characters and objects. First, a rough segmentation mask is produced by using SLIC [35] followed by hierarchical merging. This mask is then used to initialize GrabCut [36], which further refines the segmentation. The dataset also contains a clean background for each clip. Foreground characters and objects are excised, and the resulting holes are filled using PatchMatch [34]. See supplementary material for more details on the dataset.

## 5    Experiments

### 5.1    Layout Composer Evaluation

**Training.** We use the Adam optimizer (learning rate=0.001, decay factor=0.5 per epoch, weight decay=0.0001) and a batch size of 32.

**Metrics.** We evaluate layout composer using 2 metrics: (a) negative log-likelihood (NLL) of ground truth (GT) entity positions under the predicted distribution, and (b) average normalized pixel distance (coordinates normalized by image height and width) of the ground truth from the most likely predicted entity location. While NLL captures both location and scale, pixel distance only measures location accuracy. We report metrics on unseen test descriptions using ground truth locations and appearances for previous entities in the partial video.

**Feature Ablation.** The ablation study in Table 1 shows that the layout composer benefits from each of the 3 input features – text, scene context (partial video), and 2D coordinate grid. The significant drop in NLL without text features indicates the importance of entity identity, especially in predicting scale.

**Table 1. Layout Composer Analysis.** Evaluation of our model (last row) and ablations on test set. First row provides theoretically computed values assuming a uniform location distribution while making no assumptions about the scale distribution.

| Text | Scene Context | 2D Coord. Grid | Dil. Conv | NLL | Pixel Dist. |
|:---:|:---:|:---:|:---:|:---:|:---:|
| | | Uniform Distribution | | >9.704 | >0.382 |
| ✗ | ✓ | ✓ | ✓ | 9.845 | 0.180 |
| ✓ | ✗ | ✓ | ✓ | 8.167 | 0.185 |
| ✓ | ✓ | ✗ | ✓ | 8.250 | 0.287 |
| ✓ | ✓ | ✓ | ✗ | 7.780 | 0.156 |
| ✓ | ✓ | ✓ | ✓ | **7.636** | **0.148** |

The lack of spatial awareness in convolutional feature maps without the 2D co-ordinate grid causes pixel distance to approximately double. The performance drop on removing scene context is indicative of the relevance of knowing *what* entities are *where* in the scene in predicting the location of next entity. Finally, replacing vanilla convolutions by dilated convolutions increases the spatial receptive field without increasing the number of parameters improves performance, which corroborates the usefulness of scene context in layout prediction.

## 5.2   Entity Retriever Evaluation.

**Training.** We use the Adam optimizer (learning rate=0.001, decay factor=0.5 every 10 epochs) and a batch size of 30.

**Metrics.** To evaluate semantic fidelity of retrieved entities to the query caption, we measure noun, adjective, and verb recalls (@1 and @10) averaged across entities in the test set. The captions are automatically parsed to identify nouns, adjectives and verbs associated with each entity both in the query captions and target database (using GT database captions for evaluation only). Note that captions often contain limited adjective and verb information. For example, a *red hat* in the video may only be referred to as a *hat* in the caption, and *Fred standing and talking* may be described as *Fred is talking*. We also do not take synonyms (*talking-speaking*) and hypernyms (*person-woman*) into account. Thus the proposed metric underestimates performance of the entity retriever.

**Feature Ablation.** Table 2 shows that text and location features are critical to noun, adjective and verb recall. Scene context only marginally affects noun recall but causes significant drop in adjective and verb recalls.

**Effect of Auxiliary Loss.** Table 3 shows that triplet loss alone does significantly worse than in combination with auxiliary classification loss. Adding the auxiliary classification loss on either query or target embeddings improves over triplet only but is worse than using all three. Interestingly, using both auxiliary losses outperforms triplet loss with a single auxiliary loss (and triplet only) on adjective and verb recall. This strongly suggests the benefits of multi-task training in entity retrieval.

**Table 2. Entity retriever feature ablation.** Top-1 and top-10 recalls of our model (last row) and ablations while generating videos for unseen test captions.

| Query Features | | | Recall@1 | | | Recall@10 | | |
|---|---|---|---|---|---|---|---|---|
| Text | Context | Location | Noun | Adj. | Verb | Noun | Adj. | Verb |
| ✗ | ✓ | ✓ | 24.88 | 3.04 | 9.48 | 55.22 | 19.39 | 37.18 |
| ✓ | ✗ | ✓ | 60.54 | 9.5 | 11.2 | **77.71** | 39.92 | 43.58 |
| ✓ | ✓ | ✗ | 56.14 | 8.56 | 11.34 | 73.03 | 39.35 | 41.48 |
| ✓ | ✓ | ✓ | **61.19** | **12.36** | **14.77** | 75.98 | **47.72** | **46.86** |

**Table 3. Entity retriever loss ablation.** Top-1 and top-10 recalls of our model (last row) and ablations while generating videos for unseen test captions.

| Auxiliary Loss | | | Recall@1 | | | Recall@10 | | |
|---|---|---|---|---|---|---|---|---|
| Triplet | Query | Target | Noun | Adj. | Verb | Noun | Adj. | Verb |
| ✗ | ✓ | ✓ | 35.75 | 7.79 | 8.83 | 63.62 | 43.35 | 33.12 |
| ✓ | ✗ | ✓ | 51.68 | 3.8 | 8.66 | 67.86 | 25.28 | 39.46 |
| ✓ | ✓ | ✗ | 50.54 | 4.94 | 9.94 | 66.36 | 28.52 | 39.5 |
| ✓ | ✗ | ✗ | 48.59 | 3.04 | 9.34 | 65.64 | 20.15 | 37.95 |
| ✓ | ✓ | ✓ | **61.19** | **12.36** | **14.77** | **75.98** | **47.72** | **46.86** |

**Table 4. Generalization to Unseen Database Videos.** Retrieval results for CRAFT when queried against seen videos vs unseen videos.

| Video Database | Recall@1 | | | Recall@10 | | |
|---|---|---|---|---|---|---|
| | Noun | Adj. | Verb | Noun | Adj. | Verb |
| Seen (Train) | 61.19 | 12.36 | 14.77 | 75.98 | 47.72 | 46.86 |
| Unseen (Test) | 50.52 | 11.98 | 10.4 | 69.1 | 41.25 | 42.57 |

**Background retriever.** Similar to the entity recall evaluation, we computed a top-1 background recall of 57.5 for CRAFT.

**Generalization to unseen videos.** A key advantage of the embedding based text to entity video retrieval approach over text only methods is that the embedding approach can use any unseen video databases without any text annotations, potentially in entirely new domains (eg. learning from synthetic video caption datasets and applying the knowledge to generate real videos). However, this requires a model that generalizes well to unseen captions as well as unseen videos. In Table 4 we compare entity recall when using the train set (seen) videos as the target database vs using the test set (unseen) video as the target database.

**OHEM vs All Mini-Batch Triplets.** We experimented with online hard example mining (OHEM) where negative samples that most violate triplet constraints are used in the loss. All triplets achieved similar or higher top-1 noun, adjective and verb recall than OHEM when querying against seen videos ($1.8, 75.3, 8.5\%$ relative gain) and unseen videos ($1.7, 42.8, -5.0\%$ relative gain).

**Modelling Whole Video vs Entities.** A key motivation to composing a scene from entities is the combinatorial nature of complex scenes. To illustrate this

**Table 5. Human evaluation** to estimate consistency and quality of generated videos.

|  | Composition Consistency | | | Visual Quality | | |
|---|---|---|---|---|---|---|
|  | Position | Rel. Size | Interact. | FG | BG | Sharpness |
| Pixel Generation L1 | 0.69 | 0.65 | 0.55 | 0.96 | 1.44 | 1.07 |
| Ours (GT Position) | 1.69 | 1.69 | 1.34 | 1.49 | 1.65 | **2.16** |
| Ours | **1.78** | **1.86** | **1.46** | **1.98** | **1.95** | 1.82 |

point we compare CRAFT to a text-to-text based whole video retrieval baseline. For a given test caption, we return a video in the database whose caption has the highest BLEU-1 score. This approach performs much worse than our model except on verb recall (BLEU: $49.57, 5.18, 26.64$; Ours: $62.3, 21.7, 16.0$). This indicates that novel captions often do not find a match in the target database with all entities and their attributes present in the same video. However, it is more likely that each entity and attribute combination appears in some video in the database. Note that text-to-text matching also prevents extension to unseen video databases without text annotations.

### 5.3 Human Evaluation

**Metrics.** In addition to the automated recall metrics which capture semantic fidelity of the generated videos to the captions, we run a human evaluation study to estimate the *compositional consistency* of entities in the scene (given the description) and the overall *visual quality* (independent of the description). The consistency metric requires humans to rate each entity in the video on a 0-4 scale on three aspects: (a) *position* in the scene, (b) *size* relative to other entities or the background, and (c) appearance and consistency of described *interactions* with other entities in the scene. The visual quality metric measures the aesthetic and realism of the generated scenes on a 0-4 scale along three axes: (a) *foreground quality*, (b) *background quality*, and (c) *sharpness*. See supplementary material for the design of these experiments.

**Modelling Pixels vs Retrieval.** We experimented extensively with text conditioned whole video generation using models with and without adversarial losses and obtained poor results. Since generative models tend to work better on images with single entities, we swapped out the target embedding network in the entity retriever by a generator. Given the query embedding at each of the $F$ time steps, the generator produces an appearance image and a segmentation mask. The model is trained using an $L1$ loss between the masked appearance image and the masked ground truth image, and an $L1$ loss between the generated and ground truth masks. See supplementary material for more details. This baseline produced blurry results with recognizable colors and shapes for most common characters like *Fred, Wilma, Barney*, and *Betty* at best. We also tried GAN and VAE based approaches and got only slightly less blur. Table 5 shows that this model performs poorly on the Visual Quality metric compared to CRAFT. Moreover, since the visual quality of the generated previous entities affects the performance of the layout composer, this also translates into poor ratings on the composition consistency metric. Since the semantic fidelity metrics can not

**Fig. 6. Qualitative results** for CRAFT. Last row shows failures of the layout composer (left) and the entire system (right). See https://youtu.be/688Vv86n0z8 for video examples, failure cases, and visualization of predicted location and scale distributions.

be computed for this pixel generation approach, we ran a human evaluation to compare this model to ours. Humans were asked to mark nouns, adjectives and verbs in the sentence missing in the generated video. CRAFT significantly outperformed the pixel generation approach on noun, adjective, and verb recall (CRAFT $61.0, 54.5, 67.8$, L1: $37.8, 45.9, 48.1$).

**Joint vs Independent Modelling of Layout.** We compare CRAFT to a model that uses the same entity retriever but with ground truth (GT) positions. Using GT positions performed worse than CRAFT (GT:$62.2, 18.1, 12.4$; Full:$62.3, 21.7, 16.0$ Recall@1). This is also reflected in the composition consistency metric (GT:$1.69, 1.69, 1.34$; Full:$1.78, 1.89, 1.46$). This emphasizes the need to model layout composition and entity retrieval jointly. When using GT layouts, the retrieval gets conditioned on the layout but not vice versa.

# References

1. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. CoRR **abs/1312.6114** (2013)
2. Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A.C., Bengio, Y.: Generative adversarial nets. In: NIPS. (2014)
3. Bengio, Y., Mesnil, G., Dauphin, Y., Rifai, S.: Better mixing via deep representations. In: ICML. (2013)
4. Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein gan. CoRR **abs/1701.07875** (2017)
5. Radford, A., Metz, L., Chintala, S.: Unsupervised representation learning with deep convolutional generative adversarial networks. CoRR **abs/1511.06434** (2015)
6. Gregor, K., Danihelka, I., Graves, A., Rezende, D.J., Wierstra, D.: Draw: A recurrent neural network for image generation. In: ICML. (2015)
7. Reed, S.E., Akata, Z., Mohan, S., Tenka, S., Schiele, B., Lee, H.: Learning what and where to draw. In: NIPS. (2016)
8. Reed, S.E., Akata, Z., Yan, X., Logeswaran, L., Schiele, B., Lee, H.: Generative adversarial text to image synthesis. In: ICML. (2016)
9. Zhang, H., Xu, T., Li, H., Zhang, S., Huang, X., Wang, X., Metaxas, D.N.: Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. CoRR **abs/1612.03242** (2016)
10. Mansimov, E., Parisotto, E., Ba, J., Salakhutdinov, R.: Generating images from captions with attention. In: ICLR. (2016)
11. Yan, X., Yang, J., Sohn, K., Lee, H.: Attribute2image: Conditional image generation from visual attributes. In: ECCV. (2016)
12. Odena, A., Olah, C., Shlens, J.: Conditional image synthesis with auxiliary classifier gans. In: ICML. (2017)
13. Karras, T., Aila, T., Laine, S., Lehtinen, J.: Progressive growing of gans for improved quality, stability, and variation. CoRR **abs/1710.10196** (2017)
14. Hong, S., Yang, D., Choi, J., Lee, H.: Inferring semantic layout for hierarchical text-to-image synthesis. CoRR **abs/1801.05091** (2018)
15. Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. CoRR **abs/1611.07004** (2016)
16. Chen, Q., Koltun, V.: Photographic image synthesis with cascaded refinement networks. CoRR **abs/1707.09405** (2017)
17. Vondrick, C., Pirsiavash, H., Torralba, A.: Generating videos with scene dynamics. In: NIPS. (2016)
18. Villegas, R., Yang, J., Zou, Y., Sohn, S., Lin, X., Lee, H.: Learning to generate long-term future via hierarchical prediction. arXiv preprint arXiv:1704.05831 (2017)
19. Marwah, T., Mittal, G., Balasubramanian, V.N.: Attentive semantic video generation using captions. CoRR **abs/1708.05980** (2017)
20. Kwak, H., Zhang, B.T.: Generating images part by part with composite generative adversarial networks. CoRR **abs/1607.05387** (2016)
21. Zitnick, C.L., Parikh, D.: Bringing semantics into focus using visual abstraction. 2013 IEEE Conference on Computer Vision and Pattern Recognition (2013) 3009–3016
22. Zitnick, C.L., Parikh, D., Vanderwende, L.: Learning the visual interpretation of sentences. 2013 IEEE International Conference on Computer Vision (2013) 1681–1688

23. Zitnick, C.L., Vedantam, R., Parikh, D.: Adopting abstract images for semantic scene understanding. IEEE Transactions on Pattern Analysis and Machine Intelligence **38** (2016) 627–638
24. Reed, S., van den Oord, A., Kalchbrenner, N., Bapst, V., Botvinick, M., de Freitas, N.: Generating interpretable images with controllable structure. In: OpenReview.net. (2017)
25. Liu, Y., Qin, Z., Luo, Z., Wang, H.: Auto-painter: Cartoon image generation from sketch by using conditional generative adversarial networks. CoRR **abs/1705.01908** (2017)
26. Frome, A., Corrado, G.S., Shlens, J., Bengio, S., Dean, J., Mikolov, T., et al.: Devise: A deep visual-semantic embedding model. In: Advances in neural information processing systems. (2013) 2121–2129
27. Kiros, R., Salakhutdinov, R., Zemel, R.S.: Unifying visual-semantic embeddings with multimodal neural language models. arXiv preprint arXiv:1411.2539 (2014)
28. Wang, L., Li, Y., Lazebnik, S.: Learning deep structure-preserving image-text embeddings. In: Proceedings of the IEEE conference on computer vision and pattern recognition. (2016) 5005–5013
29. Faghri, F., Fleet, D.J., Kiros, J.R., Fidler, S.: Vse++: Improved visual-semantic embeddings. arXiv preprint arXiv:1707.05612 (2017)
30. Plummer, B.A., Wang, L., Cervantes, C.M., Caicedo, J.C., Hockenmaier, J., Lazebnik, S.: Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In: Computer Vision (ICCV), 2015 IEEE International Conference on, IEEE (2015) 2641–2649
31. Tsai, Y.H.H., Huang, L.K., Salakhutdinov, R.: Learning robust visual-semantic embeddings. arXiv preprint arXiv:1703.05908 (2017)
32. Caruana, R.: Multitask learning. In: Learning to learn. Springer (1998) 95–133
33. Gupta, T., Shih, K., Singh, S., Hoiem, D., Shih, K.J., Mallya, A., Di, W., Jagadeesh, V., Piramuthu, R., Shih, K., et al.: Aligned image-word representations improve inductive transfer across vision-language tasks. arXiv preprint arXiv:1704.00260 (2017)
34. Barnes, C., Shechtman, E., Finkelstein, A., Goldman, D.B.: Patchmatch: A randomized correspondence algorithm for structural image editing. ACM Transactions on Graphics-TOG **28**(3) (2009) 24
35. Achanta, R., Shaji, A., Smith, K., Lucchi, A., Fua, P., Süsstrunk, S.: Slic superpixels compared to state-of-the-art superpixel methods. IEEE transactions on pattern analysis and machine intelligence **34**(11) (2012) 2274–2282
36. Rother, C., Kolmogorov, V., Blake, A.: Grabcut: Interactive foreground extraction using iterated graph cuts. In: ACM transactions on graphics (TOG). Volume 23., ACM (2004) 309–314
37. Manning, C.D., Surdeanu, M., Bauer, J., Finkel, J., Bethard, S.J., McClosky, D.: The Stanford CoreNLP natural language processing toolkit. In: Association for Computational Linguistics (ACL) System Demonstrations. (2014) 55–60

# Supplementary Material

## Overview

The supplementary material in this PDF is organized as follows:

- Section A: Derivation of Uniform distribution inequalities in Table 1
- Section B: Details on the FLINTSTONES dataset including dataset collection and dataset diversity
- Section C: Details on the design of human evaluation tasks with Amazon Mechanical Turk interfaces

In addition, more qualitative results (in video form) can be found in the video included with this supplementary material.

## A    Derivation of Uniform distribution inequalities

The first row of Table 1 in the submitted paper provides theoretically computed values with a uniform location distribution and no assumptions about the scale distribution. Here we provide derivations for the same.

### A.1    Negative Log Likelihood

$$P(x, y, w, h) = P(x, y)P(w, h|x, y) \tag{1}$$
$$\leq P(x, y) \tag{2}$$
$$\implies \log P(x, y, w, h) \leq \log P(x, y) \text{ (Since, } P(w, h|x, y) \leq 1) \tag{3}$$
$$= \log \frac{1}{128 \times 128} \text{ (For a } 128 \times 128 \text{ image)} \tag{4}$$
$$\implies -\log P(x, y, w, h) \geq 9.704 \tag{5}$$

**Typo:** Please note a small typo in Table 1 in the main submission. The negative log likelihood for uniform distribution is $\geq 9.704$ instead of $< 9.704$.

### A.2    Pixel Distance

Let $(x_t, y_t)$ be the target (ground truth) location. The expected normalized distance (normalizing location coordinates to $[0, 1]$ range) from ground truth location given the predicted distribution is given by

$$E[\|(x - x_t)/128, (y - y_t)/128\|_2] \tag{6}$$
$$= \int_0^{128} \int_0^{128} \|x - x_t, y - y_t\|_2 P(x, y)dxdy \tag{7}$$
$$= \int_0^{128} \int_0^{128} \|(x - x_t)/128, (y - y_t)/128\|_2 \frac{1}{128 \times 128} dxdy \tag{8}$$

Note that the best case scenario is $x_t = 128/2$ and $y_t = 128/2$ (the target lies at the center of the image). Hence,

$$E[\| (x - x_t)/128, (y - y_t)/128 \|_2] \tag{9}$$

$$\geq \int_0^{128} \int_0^{128} \| x/128 - 0.5, y/128 - 0.5 \|_2 \ \frac{1}{128 \times 128} \ dx dy \tag{10}$$

$$= \int_0^1 \int_0^1 \| x' - 0.5, y' - 0.5 \|_2 \ dx' dy' \ (\text{Substituting } x' = x/128, y' = y/128) \tag{11}$$

$$= 0.382 \ (\text{Solved using Wolfram Alpha}) \tag{12}$$

## B  FLINTSTONES Dataset

### B.1  FLINTSTONES Dataset Construction

**Clip Generation.** FLINTSTONES clips are roughly three seconds (75 frames) in duration, a length chosen to capture a small number of discrete actions while limiting the number of sentences needed to describe them. Source videos for the FLINTSTONES begin as episode-length videos with no existing subdivisions. In order to assure that our clips dont span scene and shot changes, we first locate these by detecting abrupt frame-to-frame changes and subdivide between them.

**Clip Annotation.** An outline of our annotation process is shown in Figure 1. Raw clips enter the first stage of our annotation pipeline, where crowdworkers identify and localize characters by concurrently labelling them and providing their bounding boxes in three keyframes. For a small number of recurring characters (e.g. Fred, Wilma), we allow workers to select from predefined labels, while for others they write a brief description (e.g. policeman, old man in red shirt). Clips containing between 1-4 characters are passed to the next stage, where workers provide a 1-2 word description of a clip's setting (e.g. living room, park). In the third stage of the pipeline, crowdworkers write a 1-4 sentence description of the clip using the established character and setting labels. A fourth task identifies important objects mentioned in descriptions, which are annotated with bounding boxes in a fifth and final stage of the pipeline.

**Annotation Supplementation.** The dataset also provides tight segmentation masks for characters and objects, as well as clean scene backgrounds. The prohibitive cost of human-annotated masks necessitates an automated approach. First, template-matching is used to track entity positions, with the addition of a penalty term for displacement from the interpolated trajectory for stability. Tracking is run forward and backward, and the resulting trajectories are averaged to produce the final entity trajectory. Once an entities' trajectory is established, A rough segmentation mask is produced by using SLIC (Simple Linear Iterative Clustering) [35] to generate superpixels within a frame, which are

then merged hierarchically to produce regions of near uniformity. Regions overlapping an entities bounding box are joined to form a rough mask. This mask is used to initialize GrabCut [36], which further refines the segmentation. Clean backgrounds are the final component of a clip generated. Foreground characters and objects are excised, and the resulting holes are filled using PatchMatch [34]. For static backgrounds, a single median background frame is used for the entire clip, while individual frames are produced for moving backgrounds.



**Fig. 1.** An overview of FLINTSTONES Clip construction. Clip annotations are built up over several stages, each requiring and building on previous stages. Human annotation steps are denoted with a blue arrow and automated data supplementation with a green arrow.

## B.2 Dataset Diversity



**Fig. 2.** Word clouds demonstrating the diversity of the FLINTSTONES dataset. Each word cloud was created using just the top 100 words within each category, to enable a clearer visualization. (1) Top left: Top 100 Objects (2) Top right: Top 100 Characters (3) Bottom left: Top 100 verbs (4) Bottom right: Top 100 settings.

FLINTSTONES contains a wide variety of named characters, objects, and vocabulary used in describing the actions and appearances of entities and scenes as seen below. Notes: No stemming or lemmatization is performed prior to computing these statistics. Part of Speech tags were obtained using the Stanford core parser [37].

- Number of unique characters: 3897
- Number of unique objects: 2614
- Number of unique verbs: 1350
- Number of unique settings: 323

Figure 2 demonstrates a fraction of this diversity, with the most frequent character, objects, verbs and settings rendered in word clouds.

## C Design of Human Evaluation Tasks

To compute *compositional consistency* and *visual quality* (metrics mentioned in Section 5.3 in the submitted paper), we ran a human evaluation study on Amazon Mechanical Turk.

The consistency metric requires humans to rate each entity in the video on a 0-4 scale on three aspects: (a) *position* in the scene, (b) *size* relative to other entities or the background, and (c) appearance and consistency of described *interactions* with other entities in the scene.

The visual quality metric measures the aesthetic and realism of the generated scenes on a 0-4 scale along three axes: (a) *foreground quality*, (b) *background quality*, and (c) *sharpness*.

Workers were provided several examples of common defects to calibrate their ratings, in addition to written guidelines. Each video was given to three workers and their ratings averaged. To assure worker consistency between models, tasks were run simultaneously for all models and ablations. Figure 3 shows the turk interface used for *visual quality* and Figure 4 shows the mechanical turk interface used for *compositional consistency*.



**Fig. 3.** Amazon Mechanical Turk interface design used to collect the *visual quality* metric. This metric is reported in section 5.3 in the submitted paper.

**Fig. 4.** Amazon Mechanical Turk interface design used to collect the *compositional consistency* metric. This metric is reported in section 5.3 in the submitted paper.