

Explanation of the CAPGen Adversarial Patch Training Code

Overview: CAPGen and Adversarial Patch Training

This code implements **CAPGen** (Camouflaged Adversarial Pattern Generator), a method for generating adversarial patches that blend into a given environment's colors ¹ ². In traditional adversarial patch training (like AdvPatch by Thys et al. 2019), the patch's pixel values are optimized directly, often resulting in conspicuous patterns. CAPGen introduces a *palette-based parameterization*: instead of optimizing raw RGB pixels, the patch is represented as a **mix of a few base colors** extracted from the background. This ensures the patch's colors align with the environment, improving visual stealthiness ¹ ². The training pipeline uses a pre-trained YOLO object detector as the "attack target" – the patch is applied to images and optimized to **fool the detector** (reduce its detection confidence), while also maintaining realistic colors and smooth textures.

Key idea: The patch is defined by a *color probability matrix* over K base colors (from the environment). Each pixel doesn't have an independent RGB value; instead it has probabilities of being each base color. During training, these probabilities (logits) are adjusted so that each pixel leans towards one of the base colors (making the patch look camouflaged), and the patch's pattern is optimized to confuse the detector ¹ ³. Below, we break down the main components of the code and how they correspond to the CAPGen method described in the paper.

Building the Color Palette from an Image

Function: `build_palette_from_tensor(image_tensor, k)`

This function extracts a set of **base colors** from a given background image, using K-means clustering. It takes an image tensor (`[3, H, W]` in $[0,1]$) and a number of colors `k`. The steps are:

- **Reshape Image Colors:** The image's pixels are flattened into an $N \times 3$ list of RGB values (`img_np`).
- **K-Means Clustering:** It runs scikit-learn's KMeans to find `k` cluster centers in RGB space. These cluster centers are the dominant colors present in the image. For example, if `k=6`, it finds 6 representative colors from the background. This corresponds to the paper's step of *extracting specific base colors from the surrounding environment* ¹. In fact, the authors note they use clustering (K-means) to derive base colors from environment images ⁴.
- **Return Palette:** The cluster centers are returned as a PyTorch tensor of shape `[K, 3]` (each row is an RGB color). They clamp values to $[0,1]$ to ensure valid color range. This palette will remain fixed throughout training – it's a *precomputed color palette from the environment*.

In code, this is realized as:

```
img_np = image_tensor.permute(1, 2, 0).cpu().numpy().reshape(-1, 3)
kmeans = KMeans(n_clusters=k, random_state=0).fit(img_np)
centers = torch.tensor(kmeans.cluster_centers_, dtype=torch.float32) # [K,3]
centers.clamp_(0, 1)
return centers
```

So, if you provide a background image (via `config.capgen_image_path`), this function will pick out the main `k` colors from that image. These colors form the **environment palette** that the adversarial patch will use. This directly implements the paper's idea of using *base colors extracted from the environment* for the patch ¹ ⁴.

The CapGenModule Class: Differentiable Patch Generation

Class: `CapGenModule(torch.nn.Module)`

This class is the core of the CAPGen approach. It creates a *differentiable representation* of the adversarial patch using the fixed color palette. Key components and steps in `CapGenModule`:

- **Palette Buffer:** In `__init__`, it receives `palette_rgb` (the `[K,3]` tensor of base colors from above) and the patch size `(H, W)`. It registers the palette as a buffer (`self.register_buffer("palette", palette_rgb)`), meaning it's not a learned parameter but a fixed tensor on the same device. For example, if the palette has 6 colors, this is a 6×3 matrix of those RGB values.
- **Logits Parameter:** It creates a learnable tensor `self.logits` of shape `[1, H, W, K]` initialized with small random values (multiplied by 0.01). This tensor is the **color probability matrix** (logits for each pixel-color assignment). Essentially, for each pixel position in the patch ($H \times W$ positions), there are K numbers which represent that pixel's preference for each of the K palette colors. These are the parameters that the model will optimize (via gradient descent) to adjust the patch's pattern. Initially they're random (so the patch starts nearly random), but constrained to the palette's color space.
- **Forward Method:** When you call `self.capgen_mod()`, it executes the `forward()` function, which does the following:
 - **Softmax over Colors:** `r = F.softmax(self.logits / tau, dim=-1)` produces a soft probability distribution over the K colors for each pixel. This yields `r` of shape `[1, H, W, K]`, where `r[0, i, j, :]` is a K -length vector summing to 1, indicating how much pixel (i,j) uses each palette color. The division by `tau` (temperature) controls the sharpness of this distribution (explained below). At the start, with random logits, each pixel might have a roughly uniform distribution over colors. Over training, the logits will adjust so that each pixel "chooses" a particular base color or mix of colors that yields an effective adversarial pattern. This softmax corresponds exactly to the paper's **color probability matrix** optimization, where each pixel has a probability for each base color ¹. The paper denotes this matrix often as P or similar, and each pixel's color is determined by these probabilities (Eq. (3) in the paper) ⁵.

- **Mixing Colors (Convex Combination):** Next, `torch.einsum('bhwk, kc -> bhwc', r, self.palette)` computes the weighted sum of the palette colors for each pixel. In simpler terms, for each pixel location (h,w), it multiplies each palette color (an RGB vector) by the pixel's softmax weight for that color and sums them up. The result is an RGB value for that pixel. After this einsum, we get a tensor of shape `[1, H, W, 3]` (denoted as `patch` in code) – these are the actual patch pixel colors, expressed as a *convex combination of the K base colors*. This is the implementation of the paper's equation where *each patch pixel = sum_k (Probability_{k} * BaseColor_k)* ⁵. Essentially, the patch is *restricted to the span of the palette colors*.
- **Permute and Clamp:** They then permute dimensions to `[1, 3, H, W]` to make it a standard image tensor (channels-first). Finally, `patch.clamp(0,1)` ensures the pixel values stay in [0,1] range (this is mostly precautionary since the palette is in [0,1] and softmax weights sum to 1, so the combination should naturally be in [0,1] as well). The forward returns a tuple `(patch, r)`, where `patch` is the generated patch image and `r` is the soft assignment matrix (useful for debugging or optional regularization).
- **Softmax Temperature (`tau`):** The parameter `tau` (initialized e.g. to 0.07 or 0.12 by config) controls how “soft” or “hard” the color assignments are. A higher τ makes the softmax more uniform (pixels blend colors), while a low τ makes the softmax outputs peak towards one color per pixel (more one-hot). The paper explicitly uses a small τ (around 0.1) to ensure each pixel sticks to one base color ³. In this code, `CapGenModule.tau` can be updated over training (see `set_tau`) – the training loop may **anneal τ over epochs**. For example, it might start with $\tau = 0.12$ for a bit of color blending (to allow gradients from all colors), and later lower τ so that the patch becomes more discrete in color (each pixel clearly one of the base colors). This reflects the paper's strategy: “We set τ to 0.1 to keep the color of each pixel belonging to one of the base colors.” ⁶. Annealing τ is a way to gradually enforce harder assignments as training progresses.

How this works as part of a neural network: `CapGenModule` is a subclass of `torch.nn.Module`, so it integrates into the PyTorch computation graph. The learnable `self.logits` are `nn.Parameter`, which means they will accumulate gradients. In the training loop, the optimizer is set to optimize `capgen_mod.parameters()`, so it will adjust the logits based on the loss gradient. Importantly, this construction makes the patch generation **differentiable** – gradients flow from the detector's loss back into the logits via the softmax and einsum operations. This allows the patch to be optimized with gradient descent just like network weights, despite the patch ultimately being a “image” output. The advantage of this palette+logits scheme is that it constrains the patch's colors without sacrificing differentiability.

In summary, `CapGenModule` implements the method's core: maintaining a color probability matrix and producing the patch as a blend of base colors ². Each pixel's probability vector (the softmax `r`) corresponds to that pixel's membership to the palette colors. The goal (as stated in the paper) is to **optimize this probability matrix** so that the final patch both *fools the detector* and *uses only environment colors* ¹.

Training Loop: Using the Patch Module in Adversarial Optimization

The rest of the code (inside `PatchTrainer.train()`) shows how `CapGenModule` is used in an adversarial patch training routine. Here's how it works step by step, and how it aligns with the CAPGen method in the paper:

- **Initialization:** The trainer reads a config (with parameters like `patch_size`, `batch_size`, number of epochs, etc.). If `use_capgen` is True (which, by default for Option A, it is), the code expects `config.capgen_image_path` to be provided – this is the path to an image of the environment or background. It then calls `build_palette_from_tensor` on that image to get a palette of `capgen_num_colors` (for example, K=6) colors ⁴. This palette is moved to the GPU and used to create the `CapGenModule`:

```
palette = build_palette_from_tensor(bg_image, capgen_num_colors).to(device)
self.capgen_mod = CapGenModule(palette, patch_size=self.config.patch_size,
                                tau=capgen_tau_start).to(device)
```

At this point, `self.capgen_mod` contains the randomly initialized logits and the fixed palette. They even generate and save an initial patch image (`init_patch = self.capgen_mod()[0]`) for visualization – initially it will look like a random-colored pattern, but only using the palette's colors (since logits are tiny random, the softmax is nearly uniform and the patch might be an average color of the palette or a mix).

- **Optimizer and Scheduler:** If using CAPGen, the optimizer is `optim.Adam(self.capgen_mod.parameters(), lr=...)`. This means it will optimize the logits inside `CapGenModule`. (If CAPGen was off, they'd instead optimize a raw pixel tensor `adv_patch_cpu`.) A learning rate scheduler is set up as per config, but that's not specific to CAPGen.

- **Training Iterations:** For each epoch and each batch of training images:

- **Generate Current Patch:** If CAPGen is enabled, they call `adv_patch, r = self.capgen_mod()`. This produces the current adversarial patch (a `[1,3,H,W]` tensor) from the latest logits. Because the patch is generated on the fly each iteration from the logits, any gradient updates to logits from the previous iteration are now reflected in `adv_patch`. (If CAPGen were off, they'd use the raw `adv_patch_cpu` tensor, which is directly optimized).

- **Apply Patch to Images:** The code uses `PatchTransformations` and `PatchApplier` (likely to do data augmentation like random scaling, rotation, etc., and then overlay the patch onto the source images at the target object's location). Essentially, the patch is pasted on the input images (`img_batch`) according to provided target label locations (`lab_batch`). The result `p_img_batch` is a batch of images with the adversarial patch on them. They even save some sample images (like `patch_image_full.jpg`) for debugging. If weather augmentations are on, they apply brightness, rain, fog, etc., to simulate various conditions – this is to make the patch robust. Then they resize `p_img_batch` to the YOLO model's input size. At this point, we have a batch of images that contain the patch, ready to be evaluated by the detector.

- **Forward through Detector:** `output = self.darknet_model(p_img_batch)` runs the YOLO (Darknet) object detector on the patched images. The goal is to **fool this detector**. In CAPGen's context, for example, if the patch is on a person (to hide them from a person detector), the loss will be set up so that the detector's confidence for the person class is minimized. The code's `MaxDetectionScore` likely computes something like the maximum confidence of any detection of the protected object, or a related metric.

- **Compute Losses:** The code then computes several loss components:

- `det_loss = mean(max_detection)` which is the adversarial loss from YOLO's output. Typically, we want to *minimize* the detector's ability to see the object. This implementation suggests they take the detector's highest confidence and use it as a loss (so the optimizer will try to reduce this confidence). In other words, `det_loss` is high when the object is detected clearly, and by minimizing `det_loss` we force the patch to confuse the detector. This aligns with the paper's objective of maintaining adversarial effectiveness (fooling the model) ².
- `nps_loss = self.nps_calculator(adv_patch) * 0.01`. NPS stands for *Non-Printable Score* – it's a penalty used in physical adversarial patch literature to discourage colors that are not easily printable or that are very rare. It compares patch colors against a list of printable colors. Here it's weighted low (0.01) and is not a focus of CAPGen, but it helps ensure the patch colors remain in a realistic range (which in our case they already are, since they're from a real image palette).
- `tv_loss = self.tv_calculator(adv_patch) * 2.5`. TV is *Total Variation* loss – it encourages smoothness by penalizing abrupt color changes between neighboring pixels. A lower TV means the patch has smoother color transitions. They even do `torch.max(tv_loss, 0.1)` in the loss, meaning they enforce at least a minimal TV (this detail ensures the patch doesn't become completely homogeneous too soon, preventing a trivial all-one-color solution). In context, TV loss is about making the patch look more like a natural texture (less noisy), which contributes to visual stealthiness.
- **Entropy Regularizer (CAPGen-specific):** If CAPGen is used, they compute an entropy over `r` to encourage discrete color choices:

```
ent = -(r * (r + 1e-8).log()).sum(-1).mean()
ent_loss = entropy_weight * ent
```

This calculates the average entropy of the color distribution per pixel. If a pixel's distribution `r[i, j]` is spread out (high entropy), this term is larger; if `r[i, j]` is one-hot (one color has prob ~1, others ~0), the entropy is low. Multiplying by a weight `entropy_weight` (from config, could be 0.0 if not used, or some positive value) and adding to the loss means the optimizer will **penalize high entropy**. In effect, this pushes each pixel's softmax to become sharper (closer to choosing a single base color). This is exactly in line with the paper's statement: "We ensure that each pixel value corresponds to only one base color by regularizing the color probability matrix." ³. The entropy regularizer is the implementation of that idea – it's making the color probability matrix more discrete, thus each pixel will ultimately "belong" to one base color (improving camouflage by avoiding weird color blends). In the paper's terms, this improves visual harmony and makes the patch *controllable in color* ³.

- The total loss is `loss = det_loss + nps_loss + max(tv_loss, 0.1) + ent_loss`. By balancing these terms, the training aims to **fool the detector (low det_loss)** while keeping the patch printable, smooth, and close to environment colors (low nps and tv losses, plus low entropy for discrete assignment).
- **Backpropagation:** They zero the gradients, call `loss.backward()`, and `optimizer.step()`. For CAPGen, this backpropagates into `self.capgen_mod.logits`. Each iteration thus tweaks the logits so that the patch becomes a little better at evading detection and remains stealthy. Because `CapGenModule.forward()` was differentiable, the detector's gradient flows through the softmax and color mixing back to the logits. Over many iterations, this produces an optimized patch.
- **Clamping and Patch Update:** If CAPGen is *not* used (i.e. a raw pixel patch), they manually clamp `adv_patch_cpu.data.clamp_(0,1)` to keep pixel values valid. But for CAPGen, the patch is inherently clamped in `forward()`, and the palette is fixed between `[0,1]`, so they don't need an explicit clamp each time – the softmax mix will always yield a valid color. They also implement τ annealing per epoch: if `capgen_tau_end` is different from `capgen_tau_start`, they linearly interpolate τ as training progresses. For example, one might start with $\tau = 0.12$ and end with $\tau = 0.04$ after all epochs, gradually forcing more discrete color choices as the patch pattern converges. This matches the idea of keeping pixels tied to one color by the end of training ⁶.
- **Logging and Saving:** The code logs losses to TensorBoard and MLflow for monitoring. It also periodically saves the current patch image (by calling `self.capgen_mod()` to get `cur_patch`). They save images at each epoch (and even each few iterations) in `PATCH_ITERATIONS` folder and log them. This is useful to visually inspect how the patch evolves. At the end, they also save the final patch (either as an image or the raw tensor, and if CAPGen was used, they even save the trained `CapGenModule` model with `mlflow.pytorch.log_model`).

During training, as `det_loss` pushes the patch to confuse YOLO, the pattern in `logits` will form something adversarial (often a high-frequency pattern or shape). Meanwhile, the **palette constraint and entropy regularizer ensure that pattern is painted only with the environment's colors**. The result should be a patch that *looks like it belongs to the background* (camouflage) but still has an adversarial effect on the detector ¹. For example, if the background colors are greens and browns (a foliage environment), the patch will end up as a mix of those and not, say, bright random colors – making it hard for humans to spot, as intended by CAPGen ¹.

Figure: The CAPGen training pipeline (from the paper ¹ ³). A fixed set of base colors (palette) is extracted from the background. During training, a color probability matrix (soft assignments per pixel to each base color) is optimized. Each pixel's color is generated by a softmax-weighted sum of the base colors (Eq. 3 in the paper). The patch (with environment-matched colors) is then applied onto training images and passed through the object detector. The loss from the detector, combined with regularization (smoothness, etc.), updates the color probabilities. This process yields a patch that camouflages into the scene while remaining adversarial to the detector.

Relation to the Paper’s Method and Differences

The code closely follows the CAPGen method proposed in the paper ¹. Each major step in the method has a counterpart in the code:

- **Extracting Base Colors:** The use of K-Means on a background image (`build_palette_from_tensor`) corresponds to the paper’s step of extracting specific base colors from the environment ^{1 4}. This ensures the patch’s color scheme is environment-adaptive. In practice, the paper mentions using a set of images from the environment to determine base colors; the code simplifies this by using one representative image for palette extraction (which is a practical proxy for the environment). The number of colors `K` (e.g. 6) is a hyperparameter (`capgen_num_colors`) that one can experiment with, as the paper also studies how varying K affects patch performance and stealthiness ^{7 8}.
- **Color Probability Matrix (Logits) Optimization:** The core idea of optimizing a color probability matrix for the patch ¹ is exactly implemented by the `CapGenModule`’s `logits` and softmax. Each pixel’s distribution `r[i, j, :]` in the code is a row of that color probability matrix P described in the paper. The forward pass combining `r` with the palette is the realization of the paper’s Equation (3): assigning colors to each pixel by mixing base colors according to probabilities ⁵. Training this module with the detector’s loss is how the paper optimizes the patch pattern. In short, **the `CapGenModule` code is a direct translation of CAPGen’s mathematical formulation.**
- **Ensuring One Color per Pixel:** The paper emphasizes regularizing the probabilities so each pixel ends up with a single base color (for a crisp, natural look) ³. In the code, this is handled by the temperature τ (keeping it low, and possibly lowering it further) and the entropy penalty (`ent_loss`). These encourage the softmax output to approach a one-hot vector. By the end of training, each pixel of the patch should be nearly exactly one of the palette colors (making the patch essentially a mosaic of those colors, akin to camouflage). This matches the paper’s method of enforcing that constraint for visual harmony ³. If we compare, the code’s approach (entropy regularizer + tau) is one common way to do it; the paper conceptually describes it as a regularization term in the loss (with some coefficient, analogous to `entropy_weight` in code) ⁹.
- **Adversarial Objective and Other Losses:** The paper’s goal is to fool an object detector while preserving stealth. The code’s `det_loss` corresponds to attacking the detector (the loss function $\mathcal{L}(M(I+\delta(P)), \dots)$ in the paper’s formulation ¹⁰). The TV loss and NPS loss are additional considerations (not explicitly discussed in the CAPGen paper, but inherited from prior works to ensure physical realizability and smoothness). They don’t conflict with the CAPGen method; rather, they enhance the practicality: TV loss contributes to stealth by avoiding sharp artifacts, and NPS ensures the patch can be printed or realized physically (important for physical attacks like those considered in CAPGen). The paper does mention balancing *robustness*, *stealth*, and *adversarial effectiveness* with regularization coefficients ⁹ – the code’s weighted combination of losses is doing exactly that balancing act (where `det_loss` relates to adversarial effectiveness, TV/entropy to stealth, etc.).
- **Differences or Notable Implementation Details:** Overall, the code implements *Option A: Differentiable CAPGen* very much as described in the paper. One minor difference is that the paper

discusses a “fast generation strategy” where one can reuse a learned pattern and just swap in a new palette for a new environment ¹¹. The given code does not explicitly show a function to do that, but since the patch is represented by logits (pattern) and a palette, one could in principle take a trained CapGenModule, keep the logits (pattern), and change the palette to new colors – the paper claims you can quickly recolor a patch for a new background without re-optimizing the entire pattern ¹¹. This codebase seems set up to train from scratch for a given environment, so that fast transfer would require custom steps (not shown explicitly). Another small implementation detail: the code uses a single background image for palette extraction, whereas the paper might use multiple images or an averaged palette if available. But conceptually, these are the same – both get a representative set of colors.

In summary, each part of the **CAPGen paper’s method** has a clear mirror in the code: - *Base color extraction*: `build_palette_from_tensor` with K-Means ⁴. - *Color probability matrix*: `CapGenModule.logits` and softmax mixing (Eq.3) ⁵. - *Regularization for one color/pixel*: low τ and entropy loss ³. - *Adversarial training loop*: applying the patch to images and using detector loss to update the patch ².

The result of running this code is a trained adversarial patch that is **environment-adaptive** – it uses the environment’s own color palette, making it much less conspicuous to humans, yet it maintains adversarial potency against the object detector ¹. The CNN (YOLO) is “fooled” by a pattern that, to a human, might look like just another part of the scenery. This alignment between code and paper demonstrates how CAPGen is implemented in practice.

References (from code and paper): The citations in brackets (e.g., ¹) correspond to lines from the CAPGen paper or related summaries, confirming how the code reflects the described method. Each piece of code we discussed maps onto an aspect of the CAPGen technique introduced by Li *et al.* (2024) in their paper ¹ ³, ensuring that this explanation stays faithful to both the implementation and the original methodology.

¹ ³ ⁵ ⁶ ⁸ ⁹ [2412.07253] CAPGen:An Environment-Adaptive Generator of Adversarial Patches
<https://arxiv.org/pdf/2412.07253>

² ⁴ ⁷ ¹⁰ ¹¹ [Literature Review] CapGen:An Environment-Adaptive Generator of Adversarial Patches
<https://www.themoonlight.io/en/review/capgenan-environment-adaptive-generator-of-adversarial-patches>