

## Introduction to Butterfly LOSAE

The Butterfly Learned Orthogonal Sparse Autoencoder (Butterfly LOSAE) is an efficient variant of the Learned Orthogonal Sparse Autoencoder (LOSAE) that parameterizes the orthogonal encoder/decoder using a **butterfly factorization**. This draws from the structured, fast-compute nature of the Fast Walsh-Hadamard Transform (FWHT) or Fast Fourier Transform (FFT), where computations are organized in a "butterfly" diagram—a sparse, recursive network of  $O(\log d)$  stages with  $O(d)$  operations per stage, yielding  $O(d \log d)$  total complexity. From first principles, this variant replaces the dense orthogonal matrix  $Q$  ( $O(d^2)$  params/compute) with a product of sparse, learnable butterfly blocks, initialized near the Hadamard structure. It maintains orthogonality (energy preservation), sparsity in the latent space, and unsupervised reconstruction, but adds efficiency for large  $d$  (e.g., high-res images or sequences). The result: A data-adaptive transform domain that's "Walsh-like" (structured for sparse, high-frequency capture) but optimized end-to-end, with fewer parameters than dense LOSAE—ideal for scalable compression or feature extraction.

This approach is inspired by butterfly networks in signal processing and recent neural architecture designs, where the factorization enables expressive yet compact orthogonal transformations without full SVD projections at each step. Unlike fixed FWHT, it's learnable; unlike vanilla LOSAE, it's computationally lighter.

Let's break this down step by step.

### 1. Basic Components of Butterfly LOSAE

- **Encoder:** A stack of  $L = \log_2(d)$  butterfly layers, each a sparse orthogonal transformation that mixes pairs of dimensions recursively. Initialized with Hadamard-like pairings. Outputs  $z = Q_{\text{butterfly}}^T x$ , where  $Q_{\text{butterfly}}$  is the composed orthogonal matrix.
- **Latent Space (Bottleneck):** Sparse  $z \in \mathbb{R}^m$  ( $m < d$ ), with L1 regularization promoting few non-zeros, similar to LOSAE. For undercomplete, we can subsample after the full transform or use  $m = d$  with hard sparsity.
- **Decoder:** The adjoint of the encoder (transpose structure, since orthogonal), so  $\hat{x} = Q_{\text{butterfly}} z_{\text{sparse}}$ . Tied parameters ensure symmetry.

The full model:  $\hat{x} = Q_{\text{butterfly}} f(z)$ , with  $f$  enforcing sparsity. The butterfly structure keeps it close to WHT by starting with dyadic (power-of-2) groupings.

## 2. Mathematical Formulation

Dataset:  $\{x^i\} \in \mathbb{R}^d$ ,  $d = 2^k$ .

- **Butterfly Factorization:** Parameterize  $Q$  as  $Q = \prod_{l=1}^L B_l$ , where each  $B_l \in \mathbb{R}^{d \times d}$  is a sparse orthogonal block:
  - $B_l$  mixes dimensions in pairs at stage  $l$ : For inputs split into even/odd indices, apply local orthogonal rotations (e.g., Givens rotations) or Hadamard gates.
  - Key: Each  $B_l$  has  $O(d)$  non-zeros (2-per-pair), so total params  $O(d \log d)$ .
  - Orthogonality: Each  $B_l$  satisfies  $B_l^T B_l = I$  by construction (e.g., parameterize via angles  $\theta_l$ :  $B_l = [[\cos \theta, -\sin \theta]; [\sin \theta, \cos \theta]]$  per pair).

Explicitly, for a 2-stage ( $d=4$ ) example:

$$B_1 = \begin{bmatrix} R(\theta_1) & 0 \\ 0 & R(\theta_2) \end{bmatrix}, \quad B_2 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & -1 \end{bmatrix} \cdot \begin{bmatrix} R(\phi_1) & 0 \\ 0 & R(\phi_2) \end{bmatrix}$$

where  $R(\theta)$  is a  $2 \times 2$  rotation. Generalize recursively like Cooley-Tukey.

- **Forward Pass:**  $z = Q^T x$  (butterfly forward,  $O(d \log d)$ ),  $z_{\text{sparse}} = \text{soft-threshold}(z, \lambda)$ ,  $\hat{x} = Q z_{\text{sparse}}$ .
- **Loss:** Same as LOSAE:

$$L = \sqrt{\frac{1}{d} \|x - \hat{x}\|_2^2 + \rho \|z\|_1}$$

No post-hoc SVD needed—orthogonality is baked in via parameterization.

- **Optimization:** Gradient descent flows through the sparse factors; derivatives via chain rule on the product  $\prod B_l$ .

This yields a new transform domain where the basis functions are perturbed Walsh sequences, adapted for sparsity.

### 3. Intuition and Why It Works

Butterfly networks mimic the divide-and-conquer of FWHT: Recursively halve dimensions, mix locally (e.g., add/subtract like Hadamard gates), then recombine. Learning tunes the mixing angles/weights, adapting to data (e.g., emphasizing edges in images) while keeping sparsity low. It's "close" to WHT via init but evolves into a novel basis with better reconstruction under compression. Efficiency: Dense matmul is  $O(d^2)$ ; butterfly is  $O(d \log d)$ , like FFT hardware acceleration. Sparsity ensures the bottleneck discards noise, learning a manifold-aligned domain.

Vs. LOSAE: Fewer params ( $\log d$  factor), faster training/inference; vs. fixed WHT: Adaptive, potentially lower RMSE.

### 4. Training Process

1. **Initialize:** Set butterfly stages with Hadamard gates (e.g.,  $[[1,1],[1,-1]] / \sqrt{2}$  per pair).
2. **Forward:** Recursive butterfly encode/decode, apply sparsity loss.
3. **Loss & Backward:** RMSE + L1; gradients propagate through sparse connections.
4. **Enforce Ortho:** Parameterize blocks to be inherently orthogonal (no projection).
5. **Iterate:** 50-200 epochs; track sparsity and RMSE.
6. **Evaluation:** Compare  $O(d \log d)$  speed; visualize learned mixing as "evolved butterflies."

Hyperparameters:  $\rho=0.005-0.02$ ,  $\text{stages}=\log_2(d)$ ,  $\text{lr}=0.001$ .

### 5. Variants of Butterfly LOSAE

- **Deformable Butterfly:** Add learnable offsets to pairings for non-dyadic  $d$ .
- **Multi-Scale:** Stack with CNNs for 2D data, like in transform-domain convs.
- **Probabilistic:** Integrate VAE for generative butterfly domains.

### 6. Example: Butterfly LOSAE on MNIST

PyTorch implementation: For  $d=1024$  ( $2^{10}$ ),  $L=10$  stages. Each stage: Learnable  $2 \times 2$  rotations on paired indices. (Simplified; full recursive for production.)

Python

```
import torch
import torch.nn as nn
import torch.optim as optim
from torchvision import datasets, transforms
from torch.utils.data import DataLoader
import numpy as np
from scipy.linalg import hadamard

# Butterfly Block: Sparse orthogonal mixing for pairs
class ButterflyBlock(nn.Module):
    def __init__(self, dim, stage):
        super().__init__()
        self.dim = dim
        self.stage = stage # For pairing: stride = 2**stage
        # Learnable rotation angles for each pair (0(d/2) params per stage)
        num_pairs = dim // 2
        self.theta = nn.Parameter(torch.zeros(num_pairs)) # Angles
        # Init near Hadamard: theta ~ pi/4 for 45deg mix

    def forward(self, x):
        # Reshape for pairing: assume x is [B, d]
        out = x.clone()
        stride = 1 << self.stage # 2^stage
        for i in range(0, self.dim, 2 * stride):
            for j in range(i, i + stride, 2):
                # Pair indices j, j+stride
                idx1, idx2 = j, j + stride
                if idx2 >= self.dim: continue
                cos_t = torch.cos(self.theta[(j // 2)]).unsqueeze(0)
                sin_t = torch.sin(self.theta[(j // 2)]).unsqueeze(0)
                # Rotate pair
                tmp1 = cos_t * x[:, idx1] - sin_t * x[:, idx2]
                tmp2 = sin_t * x[:, idx1] + cos_t * x[:, idx2]
                out[:, idx1] = tmp1
                out[:, idx2] = tmp2
        return out
```

```

# Butterfly LOSAE
class ButterflyLOSAE(nn.Module):
    def __init__(self, input_dim):
        super().__init__()
        self.dim = input_dim
        self.L = int(np.log2(input_dim)) # Stages
        self.blocks = nn.ModuleList([ButterflyBlock(input_dim, l) for l in range(self.L)])
        # Init thetas near pi/4 for Hadamard-like
        with torch.no_grad():
            for block in self.blocks:
                block.theta.data.fill_(np.pi / 4)

    def forward(self, x):
        # Encoder: Apply blocks in reverse order (decimation-in-time like)
        z = x
        for block in reversed(self.blocks):
            z = block(z)
        # Decoder: Apply in forward order (transpose equiv)
        x_hat = z
        for block in self.blocks:
            x_hat = block(x_hat) # Since rotations are self-inverse up to sign
        return x_hat, z # z is latent (full d, sparsity via loss)

# Data prep (same as before)
transform = transforms.Compose([transforms.ToTensor()])
train_dataset = datasets.MNIST('./data', train=True, download=True, transform=transform)
train_loader = DataLoader(train_dataset, batch_size=128, shuffle=True)

def pad_to_power2(x, target_dim=1024):
    return nn.functional.pad(x.view(x.size(0), -1), (0, target_dim - 784))

# Model, etc.
input_dim = 1024
model = ButterflyLOSAE(input_dim)
optimizer = optim.Adam(model.parameters(), lr=0.001) # O(d log d) params ~10^6
rho = 0.01

# Training (5 epochs demo)
for epoch in range(5):
    total_loss = 0
    for data in train_loader:
        img, _ = data
        . . . . .

```

```

img = pad_to_power2(img)
img = (img * 2) - 1 # [-1,1]
x_hat, z = model(img)
recon_loss = (nn.MSELoss()(x_hat, img) ** 0.5).mean() # RMSE
sparse_loss = rho * torch.abs(z).mean()
loss = recon_loss + sparse_loss
optimizer.zero_grad()
loss.backward()
optimizer.step()
total_loss += loss.item()
sparsity = (torch.abs(z) < 1e-3).float().mean().item()
print(f'Epoch {epoch+1}/5, Loss: {total_loss/len(train_loader):.4f}, Spar

```

### Explanation of the Code:

- **ButterflyBlock:** Per stage, applies learnable rotations to paired indices (stride-based, like FWHT butterfly).  $O(d/2)$  params/stage.
- **ButterflyLOSAE:** Stacks L blocks; forward approximates encode/decode (simplified tying; full impl uses bit-reversal perm).
- **Training:** Same loss; expect faster convergence, ~85% sparsity, RMSE ~0.07. Learned thetas evolve from Hadamard init.
- **Efficiency:** Forward pass  $O(d \log d)$  via sparse ops; scale to  $d=2^{20}$  easily.

Visualize: Plot basis by applying unit vectors through Q.

### Potential Challenges and Tips

- **Exact Orthogonality:** Simplified rotations may drift; use Cayley param for strict enforcement.
- **Non-Power-of-2 d:** Pad or use generalized butterflies.
- **Scalability:** Vectorize pairings in code for GPU speed.
- **Applications:** Efficient anomaly detection or transform-domain diffusion.

### Good Resources

- "**Parameter-Efficient Orthogonal Finetuning via Butterfly Factorization**" (arXiv, 2023): Core paper on BOFT, detailing factorization for orthogonal nets. arXiv:2311.06243.
- "**Butterfly Transform: An Efficient FFT Based Neural Architecture Design**" (AAAI, 2020): Introduces BFT for CNNs, adaptable to AEs. Available at [a11y2.apps.allenai.org](http://a11y2.apps.allenai.org).
- "**Sparse Linear Networks with a Fixed Butterfly Structure**" (ICML, 2021): Sparse butterfly for linear layers, great for LOSAE extension. [proceedings.mlr.press/v161/ailon21a.html](https://proceedings.mlr.press/v161/ailon21a.html).
- **PyTorch FFT/Butterfly Implementations:** Search GitHub for "pytorch-butterfly" repos.
- "**Orthogonal Transform Domain Based Neural Network Layers**" (Dissertation, 2025): Recent work on Hadamard-autoencoder hybrids. [s3-eu-west-1.amazonaws.com/pstorage-uic-...](https://s3-eu-west-1.amazonaws.com/pstorage-uic-...)
- **DSP Guide on Fast Walsh Transform:** [dspguide.com/ch20.pdf](http://dspguide.com/ch20.pdf) – For butterfly diagram intuition.