# Preliminary Research Idea (11): Swin Transformer with Semantic Boundary Adherence using Swendsen-Wang Sampling

Professor Xu Yida

Machine Learning Knowledge Propagator

## Introduction and Background

Over ten years ago, I collaborated with two professors from Oxford University, Arnaud Doucet and Francois Caron, on a paper. That prior research utilized the Generalized Swendsen-Wang Algorithm for Bayesian Nonparametric Image Segmentation+. The paper is titled "Bayesian nonparametric image segmentation using a generalized Swendsen-Wang algorithm." To understand the current proposal, one must first understand Swendsen-Wang sampling+ and what it does.

## Problem Statement for Standard Swin Transformer

Standard Swin Transformer divides images into fixed $M \times M$ blocks. If an object (e.g., a cat) lies on a block boundary, its features get split. The "Shifted Window" mechanism attempts to address this by mixing cross-boundary information in subsequent layers. However, this segmentation itself remains rigid and semantically blind (i.e., lacking content-awareness).

# Intuition Behind Swendsen-Wang Sampling in Image Context

Swendsen-Wang (SW) is a Monte Carlo method used for sampling cluster configurations in Ising or Potts models. In the image context:

- Pixels are treated as "spins."

- Edges between pixels possess "bond probabilities" based on similarities in color, intensity, or learned features.

- **Clustering Step:** SW forms pixels into clusters based on these bonds. These clusters are irregular, vary in size, and naturally conform to object boundaries (similar to superpixels+).

- **Attention Mechanism:** Instead of calculating self-attention+ within fixed $7 \times 7$ blocks, it is calculated within the dynamic clusters generated by SW sampling.

- **Extension:** The method can even be applied within 3D Swin Transformers.

# Unique Advantages of the Proposed Method

If successfully implemented, this approach offers several advantages over standard Swin:

1. **Semantic Boundary Adherence:** Attention is computed within "superpixels" or image fragments, rather than arbitrary grid divisions. This reduces noise (e.g., when focusing on a foreground object, background pixels are less likely to be attended to).

2. **Long-Range Dependencies:** In standard Swin, two pixels 20 positions apart only interact after several layers of downsampling. With SW, if pixels belong to the same large, uniform region (e.g., sky), they can be immediately grouped into the same cluster, enabling instantaneous long-range communication.

3. **Adaptive Complexity:** Complex regions (high texture/high variance) can be decomposed into many small clusters, while simple regions form large clusters. This leads to a more efficient allocation of computational resources.

## Core Challenges

Combining the Swendsen-Wang (SW) algorithm with Swin Transformer presents two main challenges:

1. **Computational Complexity:** The SW algorithm involves iterative "Connected Component Labeling," which is difficult to parallelize on GPUs. This could significantly slow down training and inference compared to standard Transformers.

2. **Differentiability:** SW's clustering process involves discrete sampling operations, which obstructs backpropagation (i.e., makes it non-differentiable).