

# Why Gait Needs a Special Backbone

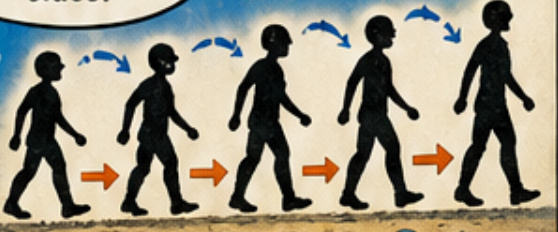
A silhouette isn't a selfie



Wait... why do CNNs keep slipping on these silhouettes?



Because gait is sparse, binary, and full of tiny motion clues!



Too flat! No grip!



### WHY IT'S HARD



Most pixels are zero. Only a few lit pixels carry the motion clues!

### GAIT TRANSFORMER BACKBONE

- ✓ END-TO-END
- ✓ GLOBAL ATTENTION
- ✓ LONG-RANGE LINKS
- ✓ CAPTURES TINY MOTION CLUES!



MOTION CLUES CAPTURED!

SMART ATTENTION!



Transformer nailed it!



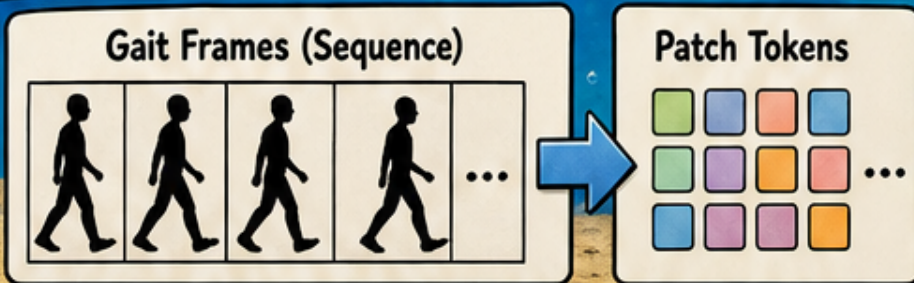
That's our special backbone!



MORE STEPS. MORE CLUES. MORE CONFIDENCE. **RECOGNIZED!**

# The Gap in Prior Transformer Gaits

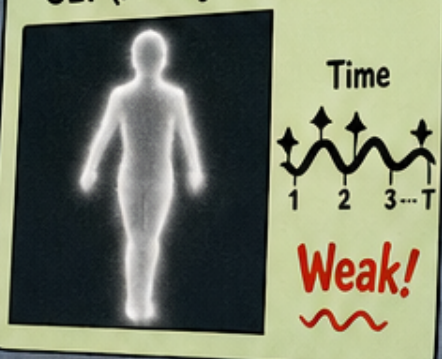
Half a solution is still wobbly



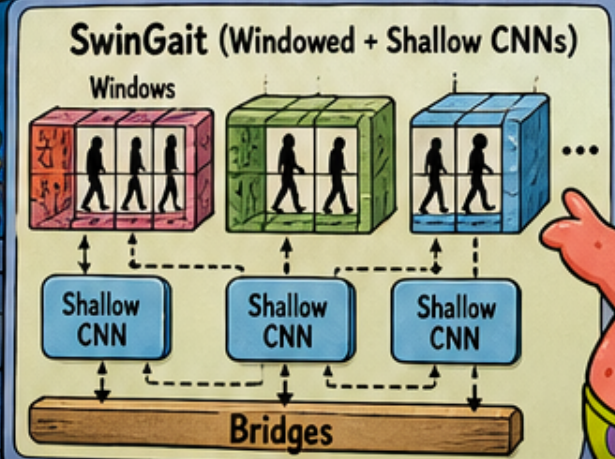
GaitViT used a vanilla ViT on GEI... but temporal modeling was weak.



GEI (Averaged Silhouette)



SwinGait helped, but leaned on shallow CNNs and bridges.



**GRAAAAH!**



The fix? An end-to-end Transformer backbone that unites all frames!



End-to-End Transformer Backbone



Global Attention Across Time

**STRONG TEMPORAL MODELING!**

End-to-end Transformers seal the gap in temporal modeling—no more wobbles!



**POP!**

FULL MODEL. FULL CONTEXT.  
**FULL CONFIDENCE!**

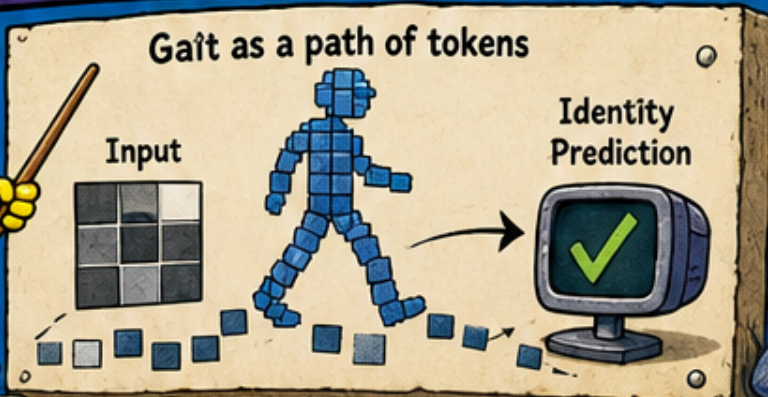
# GaT's Core Idea

One end-to-end backbone, three tricks

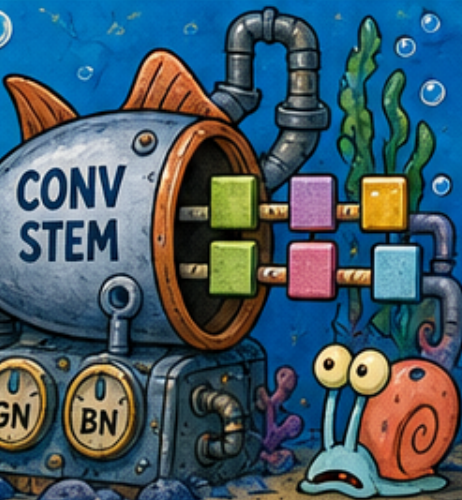
GaT as a path of tokens

Input

Identity Prediction



Hybrid patch embedding: conv stem + GroupNorm + BatchNorm.



Decomposed token mixer: short-range and long-range.

SHORT-RANGE

LONG-RANGE



WHOOOSH!

SPLASH!

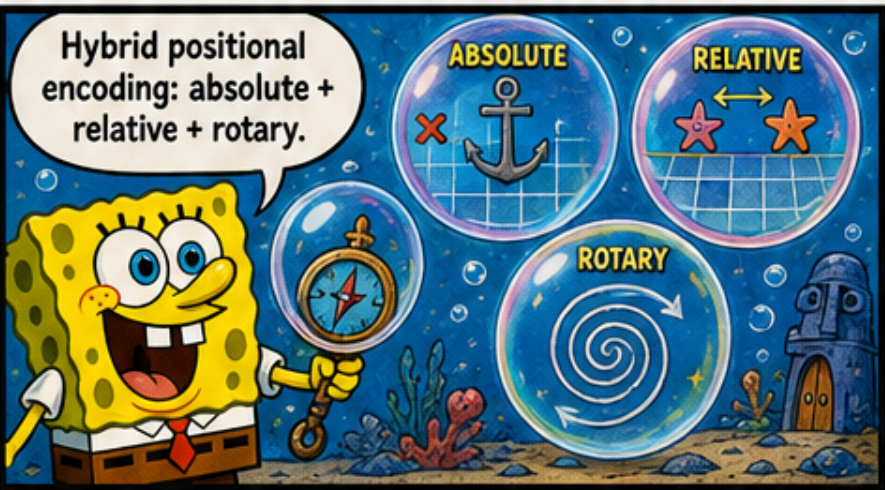


Hybrid positional encoding: absolute + relative + rotary.

ABSOLUTE

RELATIVE

ROTARY



POP!



CORAL CAFE

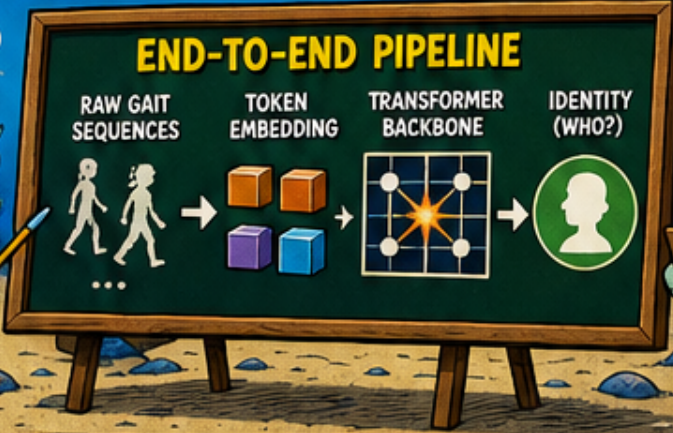
**GaT WORKS!**  
One backbone,  
Three tricks,  
Better recognition!



# How It Helps Under Data Scarcity

## Built for small, noisy gait datasets

Today I'll show how our Gait Transformer works—even with tiny, noisy data!



Sounds fancy. But we barely have data!

REEFBAY RESEARCH LAB

CORAL TOWN

### PROBLEM: DATA SCARCITY & NOISE!

Just a handful of messy walks!



### STEP 1: TURN NOISY WALKS INTO CLEAN TOKENS!

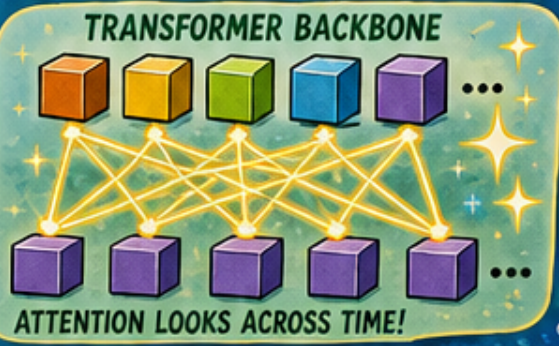
NOISY FRAMES (VERY FEW)

CLEAN TOKENS



### STEP 2: STRONG INDUCTIVE PRIORS INSIDE THE TRANSFORMER!

Strong inductive priors help when data is limited!



Underwater currents? No problem!

Different views? Bring it on!

### STEP 3: TRAIN FROM SCRATCH (NO PRETRAINING!!)

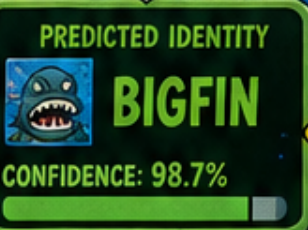
No pretraining needed—train from scratch.



Whoa, that was quick!

### RESULT: IT WORKS WITH TINY, NOISY DATA!

INPUT: FEW & NOISY



Even with so little data, I've got it!

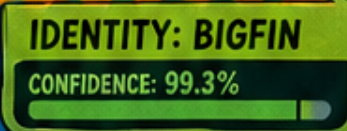
### SURPRISE TEST: WILL IT RECOGNIZE A LEGEND?

?!



IT'S BIGFIN! I KNEW IT!

**RECOGNIZED!**



THE END!

Takeaway:

# Better Gait Features, Cleaner Design

One backbone to rule the reef

Let's build the ultimate gait model!



## 1) INPUT: WALKING OVER TIME



CLEAN • UNIFIED • POWERFUL  
One backbone to rule the reef!

SPATIAL:  
How you move  
(at each moment)



+

TEMPORAL:  
How you evolve  
over time

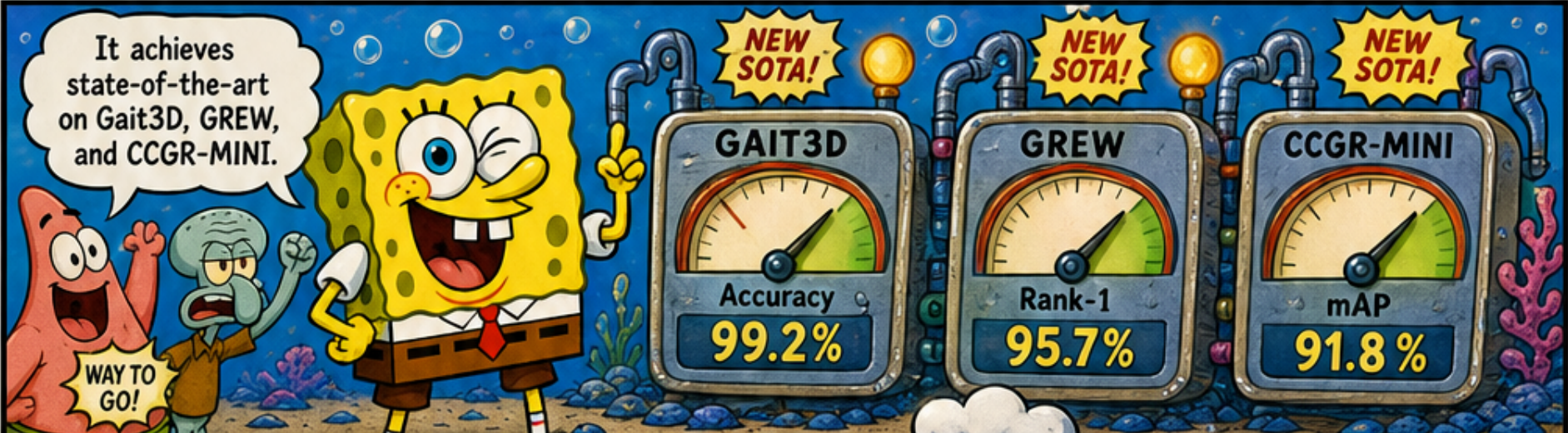


=

GaT jointly captures  
spatial and temporal  
variation.



It achieves  
state-of-the-art  
on Gait3D, GREW,  
and CCGR-MINI.



BUT THEN...

A NEW  
CHALLENGER?!

HARD  
GALI  
DATA

You thought  
it'd be that  
easy?

HA HA HA!

