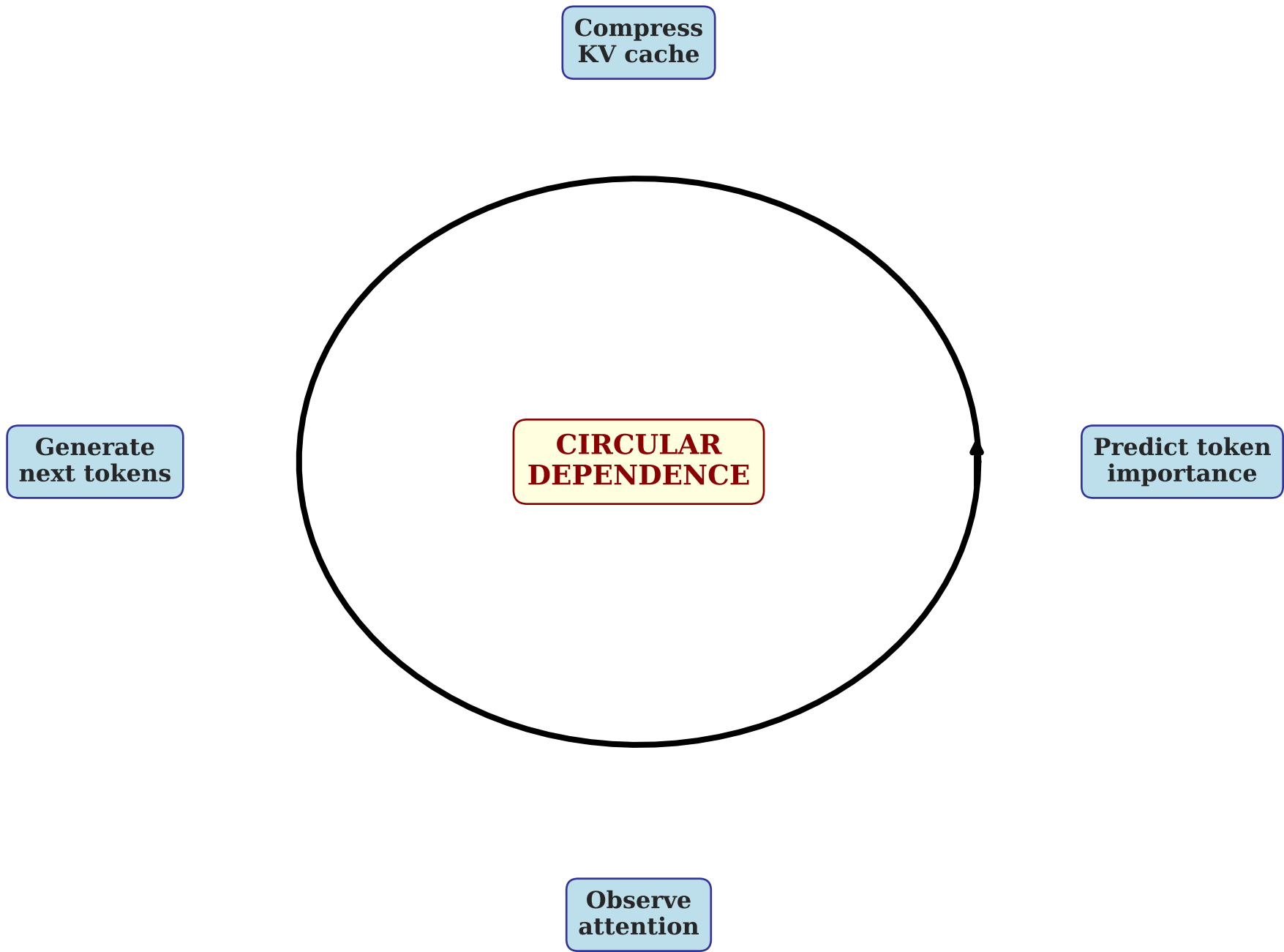# Circular Dependence in Importance Prediction

*Future importance depends on generation trajectory,*
*which depends on compression decisions*

**Compress KV cache**

**Generate next tokens**

**CIRCULAR DEPENDENCE**

**Predict token importance**

**Observe attention**

Hypothesis: Predicting which tokens will receive future attention
requires information about future queries—but those depend on tokens not yet generated.