GRABACIÓN DE CLASE

**Chard**



ruteo → fine, se rotea Al

nodo sucesor, nunca el mas cercano



todo ruteo llega Al sucesor del valor de hash de "superman"