Detailed architecture diagram:
(a) Student model pipeline
(b) Memory-aware distillation
(c) Gating network
Create using TikZ (see paper for details)

**Replace with actual figure before submission**