Avg Attention Head Flow — intervention corruption
(Top 10 heads by |marginal Δp|)