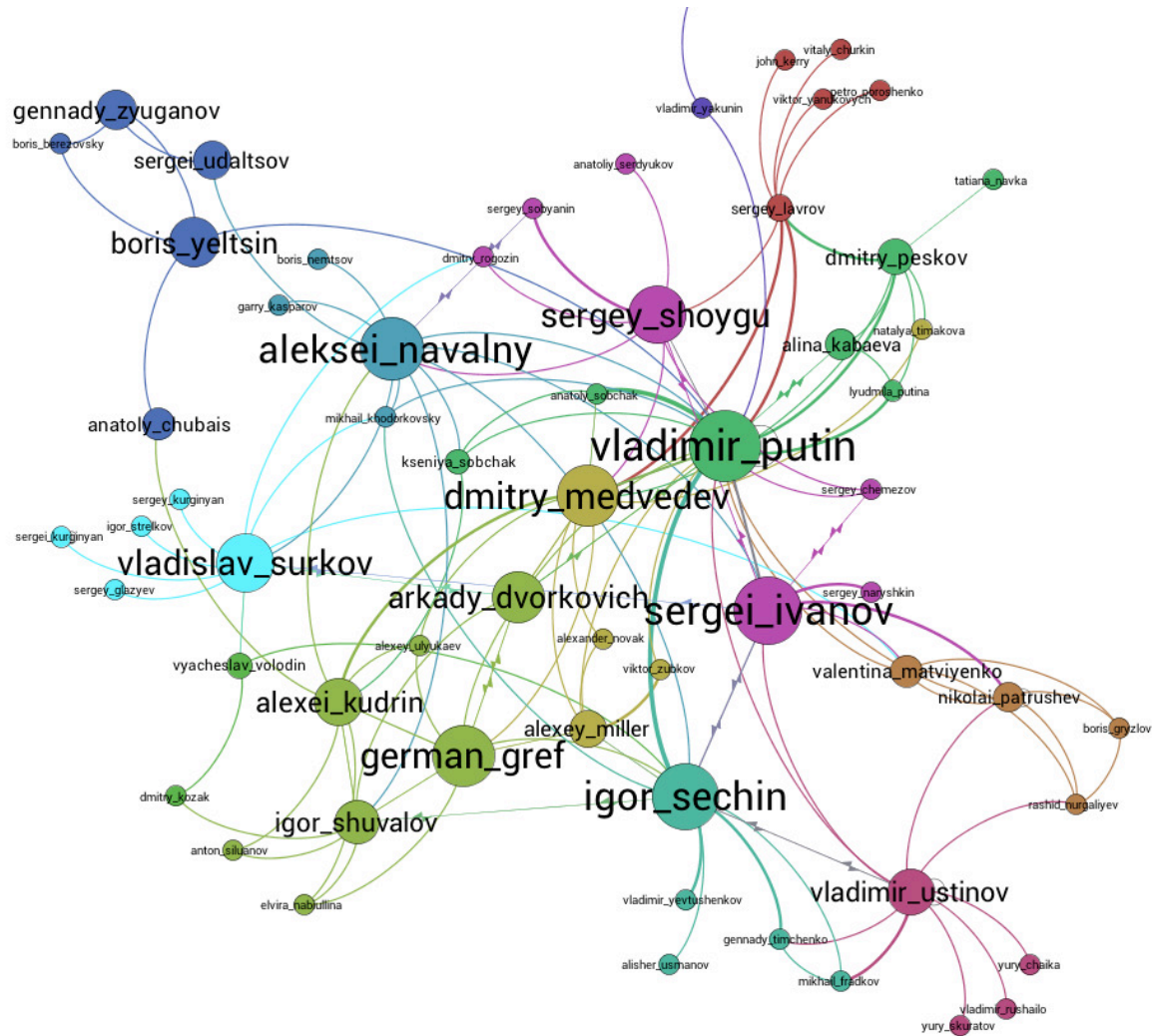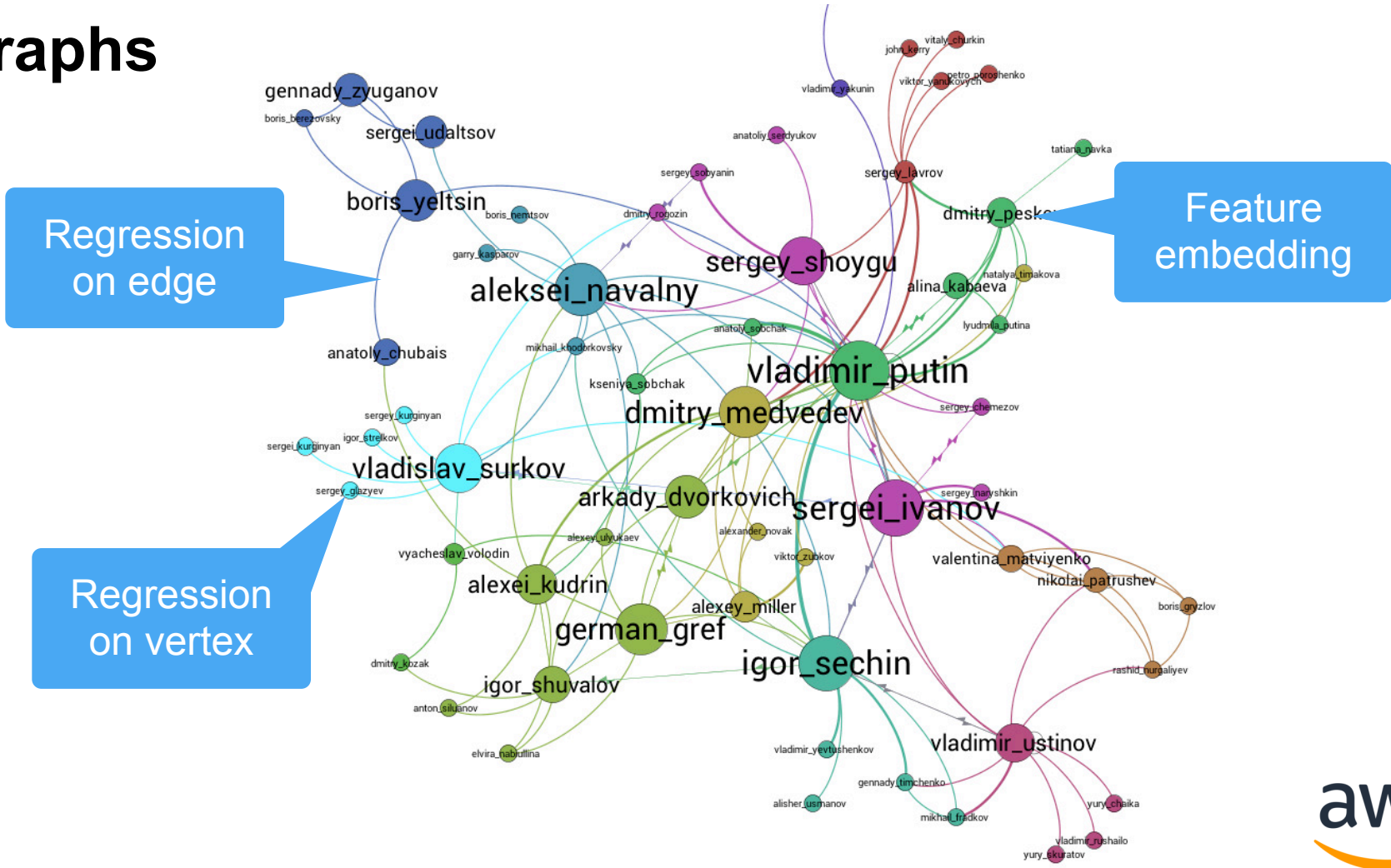# Graphs

# Beyond Sequential Dependency

Often difficult to model

- **Spatial dependency** (e.g., road network)
- **General graph dependency**
  - Spread of memes, fake news …
  - Product recommendation (users, items, vendors)
  - Relational databases
- Cannot write as sequence but needs graph. Popular choices:

**Directed** graphical model $p(x) = \prod_i p(x_i \,|\, x_{\pi(i)})$

**Undirected** graphical model $p(x) = \prod_C \psi_C(x_C)$
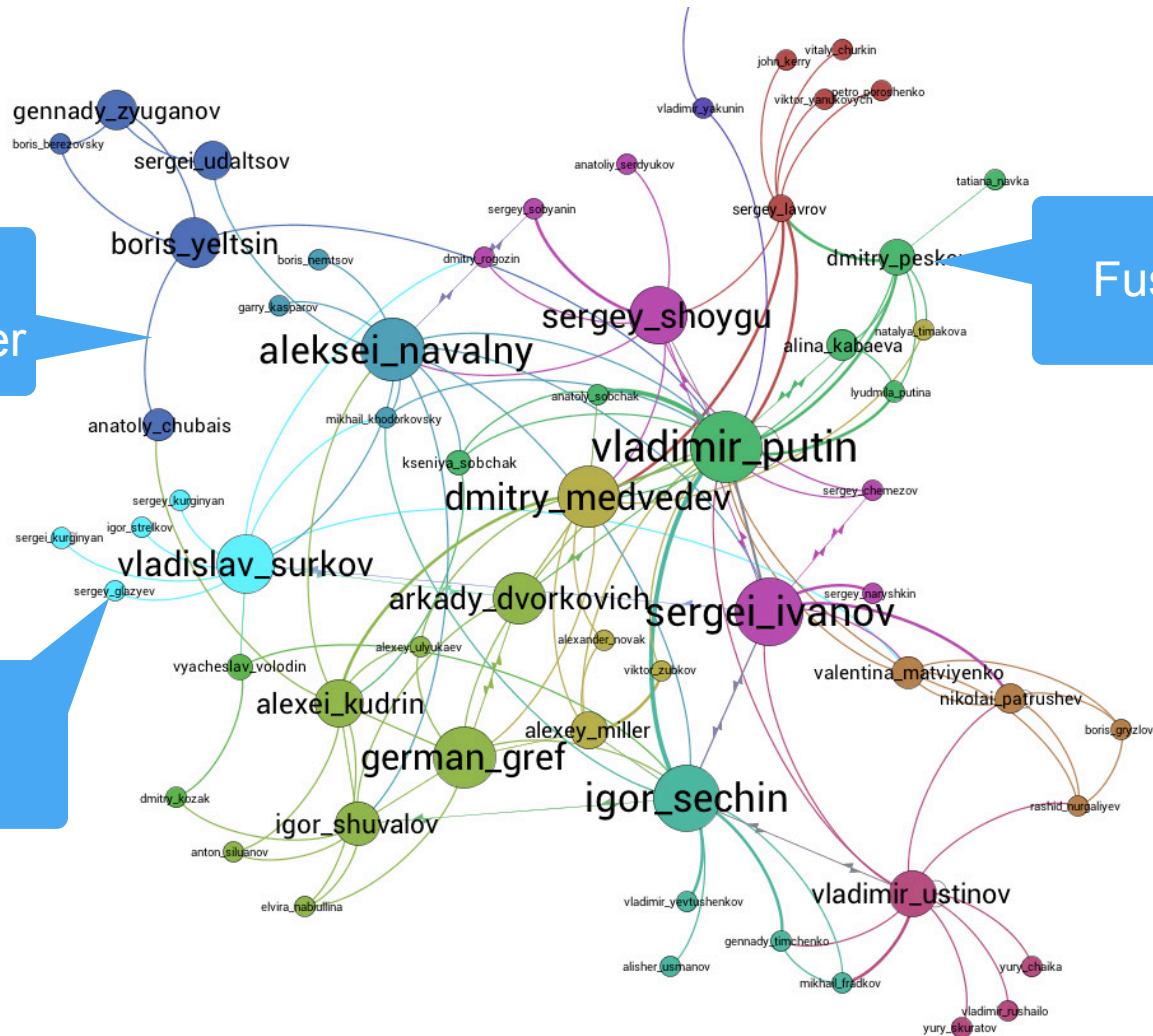
# Graphs

# Graphs



Social Recommender

Fuse data
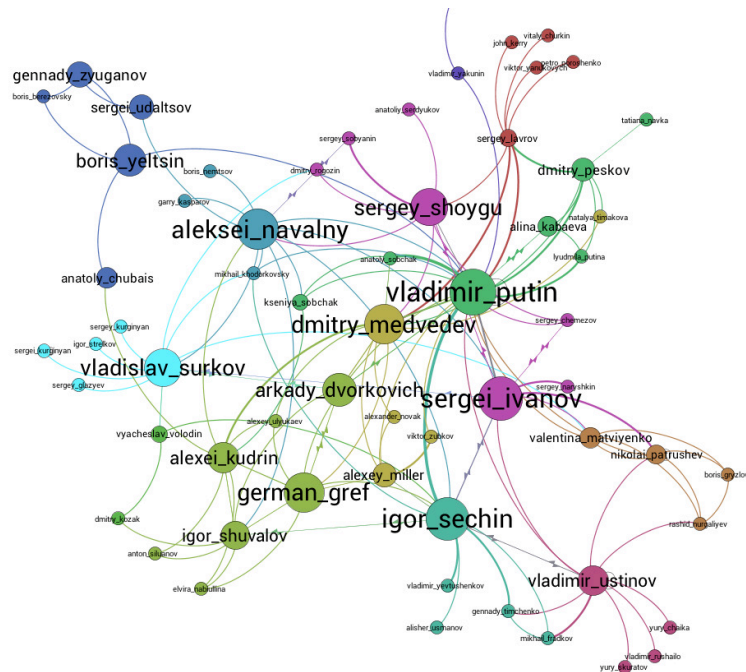
Fraud detection

# Graphs G(V,E)

- Vertices  $i \in V$  (with attributes)
- Edges  $(i,j) \in E$  (with attributes)
- Estimation problems
  - Given some vertex labels y estimate the remaining ones (e.g. fraud detection)
  - Given some edge attributes e estimate the remaining ones (e.g. link recommendation)

Vertex
Updates

# Weisfeiler-Lehman algorithm (1976)

- **Key idea**
  - Graph isomorphism … is trivial if vertices are unique
  - Make them unique by repeated hashing

$$v(i) \leftarrow h(v(i), \{v(j) \text{ with } j \sim i\})$$

  Local updates

  - Terminate once stationary
- **Machine Learning variant** (Shervashidze & Borgwardt, 2013)
  - The vertex hashes are good features $i \rightarrow \phi(i)$
  - Can prove equivalence to some graph kernels
- Crazy thought … what about vertices with attributes?

aws

# Page Rank (Page, Brin, Motwani, Kleinberg, 1990s)

- **Random surfer model with restarts**
  - Start at uniformly random location
  - Follow random link at each vertex
  - Page rank is stationary distribution
- **Self consistency equation**

$$p(i) = \frac{\epsilon}{|V|} + (1 - \epsilon) \sum_{i \sim j} \frac{p(j)}{d(j)} \iff r(i) = \epsilon + (1 - \epsilon) \sum_{i \sim j} \frac{r(j)}{d(j)}$$

- Solve by iterating fixed number of times or until convergence

Random Restart

Random Link following

Local updates

aws

# (naive) PageRank Implementation in DGL

```python
import dgl.function as fn

def pagerank_builtin(g):
    g.ndata['pv'] = g.ndata['pv'] / g.ndata['deg']
    g.update_all(message_func=fn.copy_src(src='pv', out='m'),
                 reduce_func=fn.sum(msg='m',out='m_sum'))
    g.ndata['pv'] = (1 - DAMP) / N + DAMP * g.ndata['m_sum']
```

Incoming PageRank

Random surfer

Random surfer

aws

# Belief Propagation / Message Passing

- **Graphical models** (we operate on clique graph)
- Vertex potential and directed messages

$$\phi_C = \psi_C \prod_{D \sim C} \mu_{D \to C}$$

$$\mu_{D \to C}(x_{C \cap D}) = \sum_{x_{C \setminus D}} \phi_C(x_C) \, \mu_{D \to C}^{-1}(x_{C \cap D})$$

- Finite time convergence if clique graph is a junction tree
- In practice ignore and run *for some time* on clique graph (loopy belief propagation)

aws

# Belief Propagation / Message Passing

- **Graphical models** (we operate on clique graph)
- Vertex potential and directed messages

$$\phi_C = \psi_C \prod_{D \sim C} \mu_{D \to C}$$

$$\mu_{D \to C}(x_{C \cap D}) = \sum_{x_{C \setminus D}} \phi_C(x_C) \, \mu_{D \to C}^{-1}(x_{C \cap D})$$

- Crazy thought … what if we didn't care about graphical models (they might not converge anyway)?

# This looks like vertex updates
# Can we learn them?
# Can we get vertex features?

aws

# Graph Convolutions (e.g. Kipf & Welling, 2016)

- **Basic idea**
  - Vertex (and edge) features $x_i$ and $x_{ij}$
  - Compute new vertex (and edge) features using local update function

$$v_i \leftarrow f(v_i, \{v_j, v_{ij} \text{ with } i \sim j\})$$

$$v_{ij} \leftarrow g(v_{ij}, v_i, v_j)$$

Optional (e.g. not in K&W'16)

- **Variants**
  - Run for a fixed number of steps (like graph kernel)
  - Run to convergence (like page rank)

aws

# Forward model

```python
G.ndata['feat'] = torch.eye(34)

def gcn_message(edges):
    # The argument is a batch of edges using source node's feature 'h'
    return {'msg' : edges.src['h']}

def gcn_reduce(nodes):
    # The argument is a batch of nodes with features summed over 'msg'
    return {'h' : torch.sum(nodes.mailbox['msg'], dim=1)}

class GCNLayer(nn.Module):
    def __init__(self, in_feats, out_feats):
        super(GCNLayer, self).__init__()
        self.linear = nn.Linear(in_feats, out_feats)

    def forward(self, g, inputs):
        g.ndata['h'] = inputs
        g.send(g.edges(), gcn_message)   # trigger message passing
        g.recv(g.nodes(), gcn_reduce)    # trigger aggregation at all nodes
        h = g.ndata.pop('h')             # get the result node features
        return self.linear(h)            # perform linear transformation
```

aws

# Forward Model

```python
class GCN(nn.Module):
    def __init__(self, in_feats, hidden_size, num_classes):
        super(GCN, self).__init__()
        self.gcn1 = GCNLayer(in_feats, hidden_size)
        self.gcn2 = GCNLayer(hidden_size, num_classes)

    def forward(self, g, inputs):
        h = self.gcn1(g, inputs)
        h = torch.relu(h)
        h = self.gcn2(g, h)
        return h
# First layer  – 34 inputs to 5 dimensions
# Second layer – 2 dimensional embedding for 2 groups
net = GCN(34, 5, 2)
```

aws

# Training it

```python
optimizer = torch.optim.Adam(net.parameters(), lr=0.01)
all_logits = []
for epoch in range(30):
    logits = net(G, inputs)
    # save the logits for visualization later
    all_logits.append(logits.detach())
    logp = F.log_softmax(logits, 1)
    # only compute loss for labeled nodes
    loss = F.nll_loss(logp[labeled_nodes], labels)

    optimizer.zero_grad()
    loss.backward()
    optimizer.step()
```

aws

# Graph Convolutions with state (GeniePath - Liu et al., 2018, Platanios & S, 2018)

- **Basic idea**
  - Vertex (and edge) features $x_i$ and $x_{ij}$
  - Vertex has internal state, LSTM or similar (GeniePath)
  - Compute new vertex (and edge) features using local update function

$$(v_i^{t+1}, h_i^{t+1}, c_i^{t+1}) \leftarrow \mathrm{LSTM}(v_i^t, h_i^t, c_i^t, \{v_j, v_{ij} \text{ with } i \sim j\})$$

$$v_{ij} \leftarrow g(v_{ij}, v_i, v_j)$$

Optional

  - This works better on some datasets

aws

# Graph Convolutions with state (GeniePath - Liu et al., 2018, Platanios & S, 2018)

- **Basic idea**
  - Vertex (and edge) features $x_i$ and $x_{ij}$
  - Vertex has internal state, LSTM or similar (GeniePath)
- **Minor twist (Deep Sets - Zaheer et al., 2017)**
  **All** functions on sets are nonlinear sums

$$(v_i^{t+1}, h_i^{t+1}, c_i^{t+1}) \leftarrow \text{LSTM}\left(v_i^t, h_i^t, c_i^t, g\left(\sum_{j \sim i} f(v_j, v_{ij})\right)\right)$$

$$v_{ij} \leftarrow g(v_{ij}, v_i, v_j)$$

aws

# Graph Attention Networks (Velickovic et al., 2017)

- **Basic idea**
  - Aggregation from neighboring vertices should be attention gated (self attention strategy)
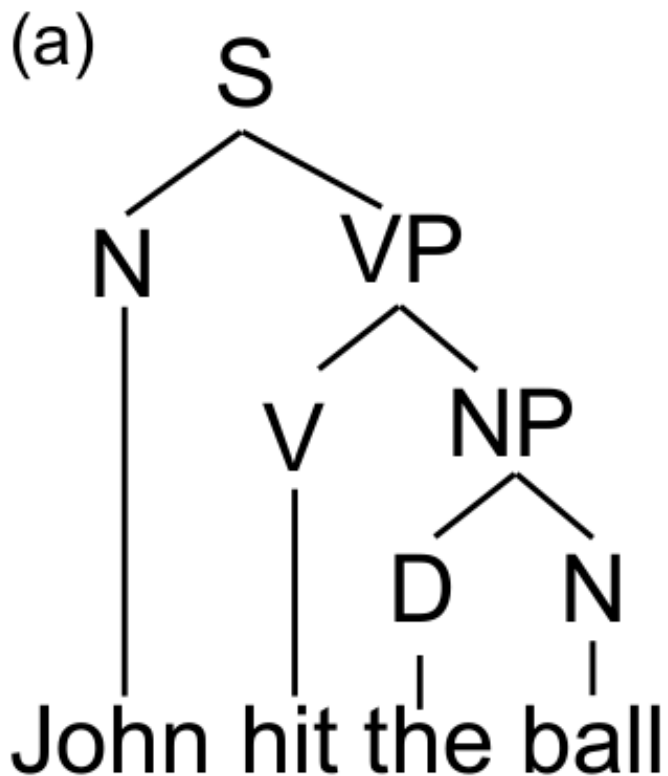  - Without attention

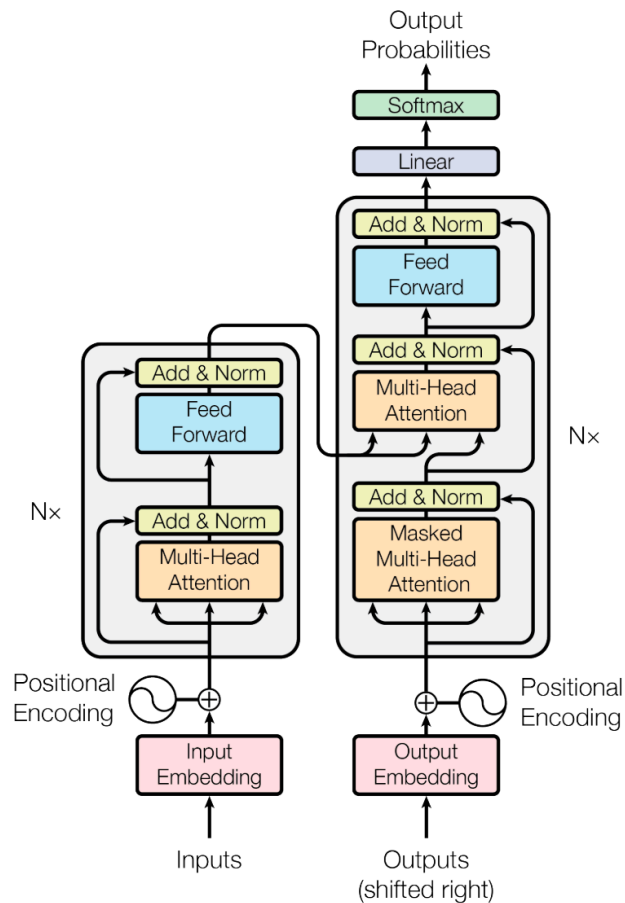$$v_i \leftarrow f\Big(v_i, \sum_{i \sim j} h(v_j)\Big)$$

  - With attention

$$v_i \leftarrow f\Big(v_i, \sum_{i \sim j} \alpha_{ij} h(v_j)\Big) \text{ where } \alpha_{ij} = \mathrm{softmax}\left(\big\{v_i^\top M v_j\big\}\right)$$
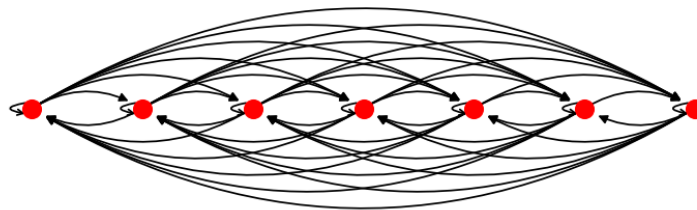
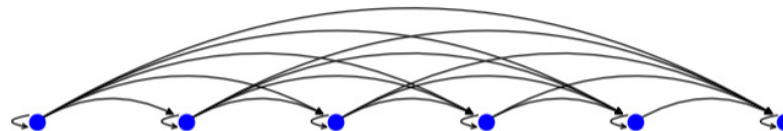  - This works better on some datasets

aws

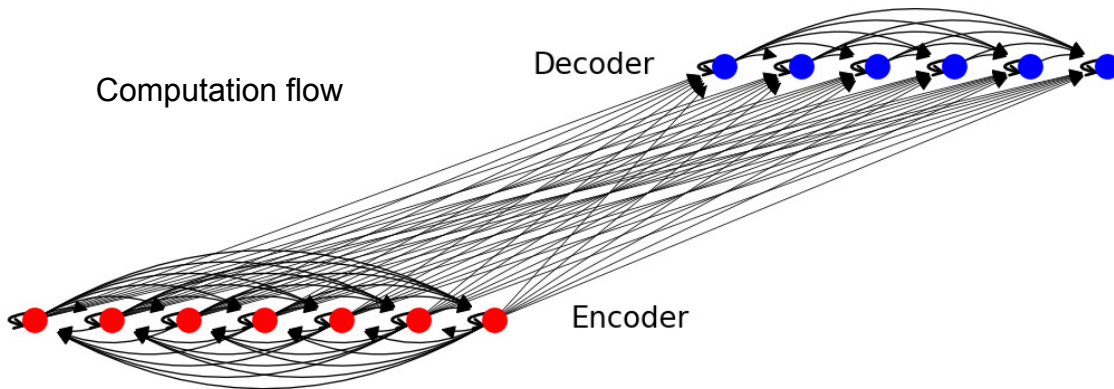# Tree-LSTM as message-passing

# Deconstruct Transformer as a graph

# Deconstruct Transformer as a graph
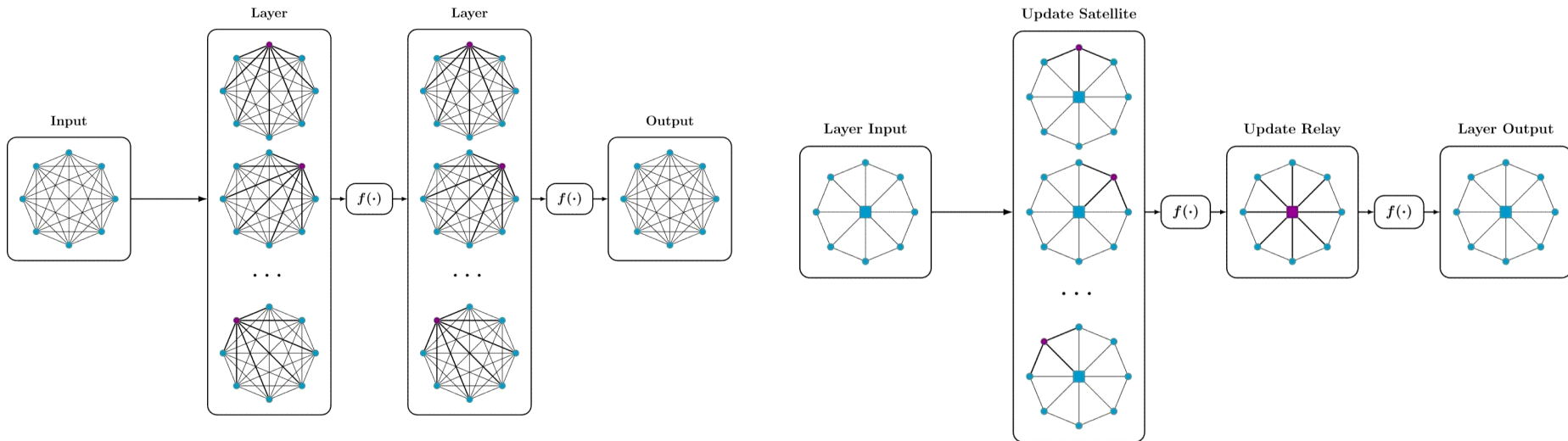


Transformer:
• data & compute hungry

Star-Transformer (in NAACL'19)
• much less data hungry
• leverages ngram prior,
• has issues with long dependency

SegTree-Transformer (in ICLR'19 RLGM)
• less data hungry
• A good compromise in between

# Making it Practical (Dai et al., 2018)

- Learning the vertex update function is **expensive**
  - Backprop over graph has to deal with entire graph quickly (6 degrees of separation kill BPTT)
  - Not much benefit in finite iterations
- Replace with fixed point iteration
  - Can learn it directly
  - Local convergence
  - **Sample vertex updates** (much smaller subset)
- This works better on some datasets

aws

# Making it Practical (Dai et al., 2018)

- Initialize
- Compute update

$$v_i \leftarrow f(v_i, \{v_j \text{ with } i \sim j\})$$

- Compute (regression) loss from embedding
- Backprop to change
  - f via loss
  - f via self consistency