

Manifold-Constrained Hyper-Connections (mHC): Architecture and Theory

Manifold-Constrained Hyper-Connections (mHC) is a recent extension of the standard Transformer residual block designed to widen the residual path while preserving the stability of identity mappings. In a conventional Transformer, each layer computes

$$x_\ell = F(x_{\ell-1}) + x_{\ell-1},$$

so that the input $x_{\ell-1}$ is passed directly to the output. This *identity mapping* ensures stable training in very deep networks ¹ ². In contrast, **Hyper-Connections (HC)** [Zhu et al. 2024] expand the residual stream into k parallel channels, introducing learnable weights to mix information across these channels. Specifically, HC replaces the simple skip-connection by three learnable matrices W_r , W_{out} , and W_{in} : $W_r \in \mathbb{R}^{k \times k}$ mixes the k residual streams, $W_{\text{out}} \in \mathbb{R}^k$ collapses them back into a d -dimensional output, and $W_{\text{in}} \in \mathbb{R}^{d \times k}$ writes the output back into the streams for the next layer ³. This greatly increases the *information bandwidth* of the residual path, but it also breaks the strict identity property, leading to instability at scale ¹ ³.

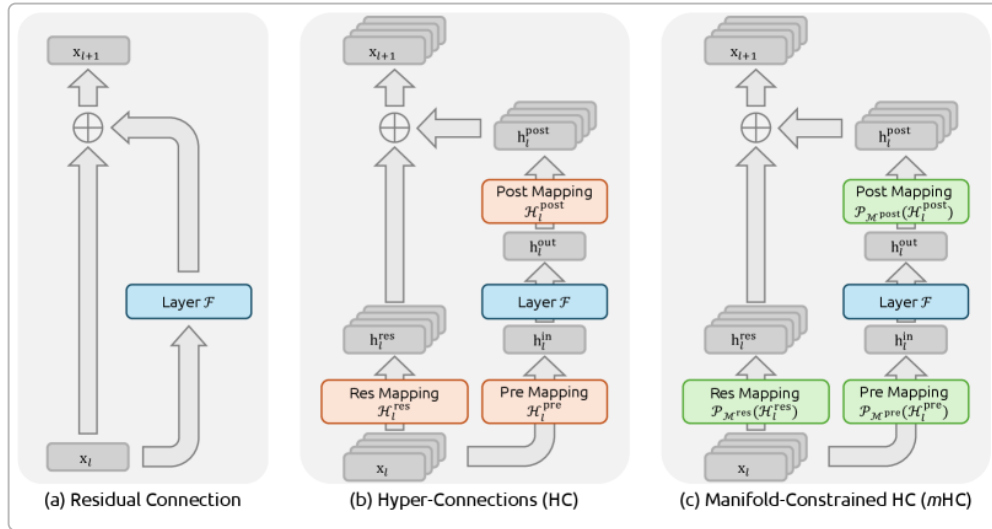


Figure: Illustration of residual paradigms. (a) Standard residual block with identity skip. (b) Hyper-Connections (HC) block with k -channel residual streams (orange) and learnable mix weights W_r , W_{out} , W_{in} . (c) Manifold-Constrained Hyper-Connections (mHC) enforces W_r to lie on the doubly-stochastic manifold (green), restoring the mean-preserving identity property while keeping a widened residual stream ¹ ³.

In mHC, the **core idea** is to constrain the HC mixing matrix W_r to the *Birkhoff polytope* of doubly stochastic matrices. Concretely, each $W_r \in \mathbb{R}^{k \times k}$ is enforced to satisfy

$$W_r \mathbf{1} = \mathbf{1}, \quad W_r^T \mathbf{1} = \mathbf{1}, \quad W_{r,ij} \geq 0,$$

so that every row and column sums to 1 ⁴. This constraint makes W_r a convex combination of permutation matrices. Equivalently, multiplying by W_r performs a weighted averaging of the k streams. By construction, this **conserves the feature mean** (the average across streams) and bounds the spectral norm to at most 1 ¹ ². In other words, each stream's signal is non-expansive and the *composition* of such mappings across layers remains identity-like. Indeed, because the set of doubly-stochastic matrices is *closed under multiplication*, the product $W_r^{(\ell_1)} W_r^{(\ell_2)} \dots$ stays doubly-stochastic, preserving stability throughout the network depth ¹ ².

Thus, **mHC restores the identity mapping property** by turning the multi-stream residual connection into a *stable feature fusion* layer. The nonnegative, doubly-stochastic W_r acts as a bounded mixing: no signal can explode or vanish across many layers (its spectral norm is ≤ 1) ¹. Meanwhile, the residual stream is still k -times wider, so the model can learn richer cross-layer interactions. The overall architecture of an mHC-enhanced Transformer layer is identical to the standard Transformer except that every skip-connection uses this k -stream mechanism. In each layer ℓ , one maintains a hidden state $H^{(\ell)} \in \mathbb{R}^{k \times d}$ (the residual streams). The layer input $x_{\ell-1} \in \mathbb{R}^d$ is expanded into k copies to form H_{in} , and the output of the layer's feedforward/attention block is added back into the streams. The update within one block is therefore: 1. **Mix streams:** $H' = W_r H_{\text{in}}$, where W_r is projected to be doubly-stochastic ¹.

2. **Collapse to output:** $x_\ell = x_{\ell-1} + W_{\text{out}} H'$, where W_{out} (a nonnegative $1 \times k$ vector) sums the streams into the d -dimensional output ³.

3. **Write back:** Update the streams by $H_{\text{out}} = H' + W_{\text{in}} y$, where $y = F(x_{\ell-1})$ is the layer's main output and W_{in} (a $k \times d$ nonnegative matrix) projects that output into the streams ³.

In summary, an mHC block behaves like a normal residual block except that its skip-connection is replaced by a **manifold-constrained hyper-connection** involving W_r , W_{out} , W_{in} . A concrete pseudocode of the forward pass is:

```
def mHC_block(x_prev, H_prev):
    # x_prev: d-dimensional input; H_prev: kxd residual streams
    # 1) Flatten streams to form context vector (for dynamic mapping) and
    compute W_r:
        a_dyn = W1 @ H_prev.flatten() + b1
        a_stat = W2 @ H_prev.flatten() + b2
        W_hat = sigmoid(a_dyn + a_stat)          # ensure positivity
        W_exp = exp(W_hat)                      # make entries positive
        # 2) Sinkhorn-Knopp projection to make W_exp doubly-stochastic:
        W = W_exp
        for t in range(N_iter):
            W = W / W.sum(dim=1, keepdims=True) # normalize rows
            W = W / W.sum(dim=0, keepdims=True) # normalize columns
        # 3) Apply residual mixing:
        H_mixed = W @ H_prev                    # k*k times kxd -> kxd
        # 4) Compute main layer output and combine:
        y = F(x_prev)                          # e.g. attention or FFN (dimension
    d)
        x_new = x_prev + W_out @ H_mixed        # collapse streams to d-dim and add
    skip
```

```
# 5) Update streams:
H_new = H_mixed + (W_in @ y).reshape(k,d)
return x_new, H_new
```

Here **Step 2** implements the **Sinkhorn-Knopp** algorithm ⁵ to project the unconstrained matrix $\exp(W_{\text{hat}})$ onto the doubly-stochastic manifold. Each iteration alternately normalizes all rows and all columns to sum to 1. After enough iterations (the paper uses $N_{\text{iter}} = 20$), W is (approximately) in the Birkhoff polytope ⁵. In training, this projection is part of the forward pass; gradients flow through the Sinkhorn iterations (the authors implement a custom backward pass) so that the original, un-projected weights are learned end-to-end. No extra loss term is needed – the constraint is enforced directly by the projection.

Mathematical Properties and Formulation

Formally, letting $\mathcal{B}_k = \{W \in \mathbb{R}^{k \times k} : W\mathbf{1} = \mathbf{1}, W^T\mathbf{1} = \mathbf{1}, W_{ij} \geq 0\}$ be the Birkhoff polytope, mHC enforces $W_r^{(\ell)} \in \mathcal{B}_k$ for every layer ℓ ⁴. Equivalently, each layer’s mixed output is a convex combination of its k input streams. This yields three key benefits ¹ ²:

- **Spectral-Norm Bound:** Any doubly-stochastic matrix has spectral norm ≤ 1 ⁴. Thus W_r is a non-expansive map. The hidden signal cannot amplify (nor vanish) exponentially across layers, preventing exploding/vanishing gradients.
- **Identity Closure:** Since \mathcal{B}_k is closed under multiplication, the product of successive W_r matrices remains doubly-stochastic ². Therefore the *global* effect of many layers still acts like a weighted identity mapping on average, preserving the feature mean throughout the network ¹ ².
- **Geometric Interpretation:** By Birkhoff-von Neumann, a doubly-stochastic W is a convex mixture of permutation matrices ⁴. mHC thus continuously mixes streams in a balanced way. Over many layers, information from all streams gets gradually fused monotonically.

Importantly, the input-output identity is recovered as a special case: if $k = 1$, the only “matrix” is $[1]$ which exactly matches the original skip-connection. In general $k > 1$ allows richer mixing, but the conservation of row/column sums ensures that “on average” the identity path remains intact ¹. The input and output mappings W_{in} and W_{out} are also constrained to be nonnegative, avoiding any cancellation across streams. These combined manifold constraints guarantee that mHC behaves like a *well-conditioned* generalization of residual connection.

Parameterization and Learning

In practice, each weight matrix is implemented using a *dynamic + static* decomposition, as in the original Hyper-Connections formulation ⁶. Concretely, at layer ℓ we flatten the hidden streams $H^{(\ell-1)}$ into a vector and compute two sets of coefficients via small linear projections: a dynamic component (input-dependent) and a static bias. A sigmoid nonlinearity then enforces positivity. For example, one sets

$$W_{\text{hat}} = \sigma(W_a \text{vec}(H) + W_b \text{vec}(H) + b),$$

where σ is sigmoid and b is a learnable bias. Learnable **gating scalars** α, β (initialized small) can weight the dynamic vs. static parts ⁶. This yields an unconstrained positive matrix $W_{\text{exp}} = \exp(W_{\text{hat}})$. We then perform the Sinkhorn normalization: for iteration $t = 0, \dots, n - 1$,

$$W^{(t+1)} = C(R(W^{(t)})), \quad R(M)_{ij} = \frac{M_{ij}}{\sum_j M_{ij}}, \quad C(M)_{ij} = \frac{M_{ij}}{\sum_i M_{ij}},$$

alternating row and column normalization. As $t \rightarrow \infty$, $W^{(t)} \rightarrow W_r$ becomes doubly-stochastic ⁵. In mHC, the authors fix $n = 20$ to get an approximate projection ⁵. This W_r is then used to mix streams. Throughout training, gradients are back-propagated through all these steps (the Sinkhorn steps are differentiable with a custom backward pass) so that the underlying parameters W_a, W_b, b are learned to produce stable W_r .

The other mappings $W_{\text{out}} \in \mathbb{R}^{1 \times k}$ and $W_{\text{in}} \in \mathbb{R}^{k \times 1}$ are learned similarly but without a global normalization: typically one ensures these are nonnegative (e.g. by ReLU or Sigmoid) so that adding and writing into streams does not cancel signals. In ablations, the authors found W_r (the residual mixing) was the most important of the three mappings ⁶. No additional loss terms are used: the model is trained end-to-end on the usual language-modeling objective (e.g. cross-entropy over tokens). The manifold constraint itself is enforced solely by the above parameterization and Sinkhorn projection.

Training Stability and Efficiency

Empirically, mHC enables stable training of very large models. In large-scale experiments (up to 27B parameters), unconstrained HC exhibited catastrophic loss spikes and gradient explosions. In contrast, mHC-trained models show smooth loss curves and bounded gradient norms ⁷ ¹. In other words, the manifold constraint effectively restores the “identity-map conservation” of residual nets. The loss gap between HC and mHC grows as models scale, confirming that mHC removes a structural instability ⁷. Across multiple benchmarks, mHC models consistently outperformed standard Transformers of the same size, indicating the gains are architectural, not just hyperparameter tweaks ¹ ⁷.

Of course, widening the residual stream increases memory I/O. To make mHC practical, DeepSeek engineers applied several optimizations ⁸ ⁹. Key strategies include **kernel fusion** (combining the normalization, projection and residual merge into efficient GPU kernels), **selective recomputation** (discarding and later recomputing intermediate activations for the streams), and overlapping communication in pipeline parallelism. These optimizations limit the added overhead: in reported runs a 4× wider stream incurred only ~6–7% extra training time ⁸. The result is that mHC can be deployed at scale with reasonable cost, unlocking a higher-capacity residual topology without the usual training collapse.

Related Concepts

- **Residual and Highway Networks:** mHC can be seen as a generalization of Highway/Residual designs ³. Highway networks introduced gated skip connections, and later designs like DenseNet or Deep Layer Aggregation added cross-layer links. HC/mHC instead add *wide, weighted* skip connections.
- **Hypernetworks:** Though similarly named, traditional *Hypernetworks* (Ha et al. 2016) are different: they use one neural network to **generate weights** for another ¹⁰. mHC’s “hyper-connections” do not generate other weights; rather, they **mix** the residual streams with trainable matrices.
- **Manifold Constraints in Learning:** Constraining weights to manifolds is a known technique. For instance, enforcing orthonormal weight matrices (on the Stiefel manifold) ensures spectral norm = 1 in layers, which also improves stability. mHC’s use of the *doubly stochastic* manifold is a novel

instance of this principle ¹ ¹¹. By projecting onto a Riemannian manifold with known closure properties, mHC provides rigorous guarantees about training dynamics.

- **Related Works:** Beyond HC, other “multi-stream” architectures have been proposed (e.g. Residual Matrix Transformer, MUDDFormer, etc. cited in ¹² ¹¹), but most lacked a manifold constraint and therefore hit scalability issues. The mHC-GNN work ¹¹ ² adapts these ideas to graph networks, and theoretical analyses there mirror mHC’s claim: constraining to a doubly-stochastic mix *prevents feature collapse* in deep message passing, just as it prevents gradient explosion here.

End-to-End Walkthrough

Putting it all together, the *end-to-end* forward pass of an mHC-enhanced Transformer block is:

1. **Pre-Normalize** the input $x_{\ell-1}$ (as in a standard Transformer).
2. **Compute Residual Mixing:** Expand $x_{\ell-1}$ into k streams to form a hidden matrix $H^{(\ell-1)}$. Compute the combined coefficients (dynamic + static) and apply the Sinkhorn projection to get $W_r^{(\ell)} \in \mathcal{B}_k$ ¹³ ¹.
3. **Mix Streams:** Multiply $H^{(\ell-1)}$ by $W_r^{(\ell)}$ to get $H_{\text{mixed}}^{(\ell)}$.
4. **Feedforward/Attention:** Process $x_{\ell-1}$ through the attention or feedforward sublayer to get y_ℓ .
5. **Merge Outputs:** Collapse the mixed streams via $W_{\text{out}}^{(\ell)}$ and add to the input: $x_\ell = x_{\ell-1} + W_{\text{out}}^{(\ell)} H_{\text{mixed}}^{(\ell)}$.
6. **Update Streams:** Write the new output into the streams via $W_{\text{in}}^{(\ell)}$: set $H_{\text{new}}^{(\ell)} = H_{\text{mixed}}^{(\ell)} + W_{\text{in}}^{(\ell)} y_\ell$.
7. **Proceed:** Pass x_ℓ to the next layer (and $H_{\text{new}}^{(\ell)}$ as the updated streams).

Throughout, all weight parameters $\{W_a, W_b, b, \alpha, \beta\}$ that define W_r , as well as $W_{\text{in}}, W_{\text{out}}$, are learned by gradient descent on the usual language-model objective. The Sinkhorn projection is simply part of the forward graph. In this way, mHC ties together theory and practice: the **theoretical requirement** (rows/columns sum to 1) is enforced algorithmically at each step, and the **implementation** (kernels, fused operations) ensures it runs efficiently.

In summary, Manifold-Constrained Hyper-Connections widen the residual pathway (like HC) but **rigidly constrain** the connections to a stability-preserving manifold (via Sinkhorn) ¹ ¹¹. This yields a new class of Transformer layers that combine high capacity with provable training stability, suggesting a promising path for next-generation deep networks.

Sources: The above explanation is based on the *mHC* paper by Xie et al. (2026) ¹ ¹³ and related references ³ ¹¹ ¹⁰. All equations and figures are adapted from these sources.

¹ ³ ⁴ ⁵ ⁶ ⁷ ⁸ ⁹ ¹² ¹³ [2512.24880] mHC: Manifold-Constrained Hyper-Connections
<https://arxiv.org/html/2512.24880>

² ¹¹ mHC-GNN: Manifold-Constrained Hyper-Connections for Graph Neural Networks
<https://arxiv.org/html/2601.02451v1>

¹⁰ [1609.09106] HyperNetworks
<https://arxiv.org/abs/1609.09106>