

377 A Description of architectures

378 **Motion Compensation** We compare our method to the traditional motion-compensated coding
379 approach that forms the core of inter-picture coding in well established compression standards such
380 as MPEG. Block matching is an essential component of these standards, allowing the compression of
381 video content by up to three orders of magnitude with moderate loss of information. For each block
382 in a frame, typical coders search for the most similar spatially displaced block in the previous frame
383 (typically measured with MSE), and communicate the displacement coordinates to allow prediction
384 of frame content by translating blocks of the (already transmitted) previous frame. We implemented
385 a “diamond search” algorithm [29] operating on blocks of 8×8 pixels, with a maximal search
386 distance of 8 pixels which balances accuracy of motion estimates and speed of estimation (the search
387 step is computationally intensive). We use the estimated displacements to perform causal motion
388 compensation (cMC), using displacement vectors estimated from the previous two observed frames
389 (x_{t-1} and x_t) to predict the *next* frame (x_{t+1}) rather than the current one (as in MPEG).

390 **Complex Steerable Pyramid** We consider a fixed multiscale oriented representation of image
391 content: a steerable pyramid [11, 10] covering 16 orientations and 5 scales on the DAVIS dataset
392 (resp. 16 orientations and 4 scales on VanHateren dataset). This choice of number of orientations and
393 number of scales maximizes prediction performance on the corresponding datasets.

394 **Polar Predictor** We use 16 pairs of convolutional channels with filters of size 17×17 pixels,
395 without biases (no additive constants). For the multiscale version, the representation is computed
396 inside a fixed Laplacian pyramid [9]. We used 4 scales for the DAVIS dataset (and respectively 4
397 scales for the VANH dataset). Within this multiscale representation, the learned filters are applied
398 with zero padding (ie. "same" boundary condition).

399 **Quadratic Predictor** We consider a multiscale architecture (cf. Polar Predictor) with 16 groups of
400 4 convolutional filters for the analysis (f_w) and synthesis (g_w) mappings. The quadratic predictor (p_w)
401 operates on groups of 4 coefficients and contains 12 quadratic units. It is therefore more expressive
402 than the Polar Predictor architecture and contains phase advance as a special case.

403 **Vanilla CNN** Finally, we implemented a more direct convolutional neural network predictor (CNN),
404 that maps two successive observed frames to an estimate of the next frame [12]. For this, we used a
405 CNN composed of 20 stages, each consisting of 64 channels, and computed with 3×3 filters without
406 additive constants, followed by half-wave rectification. To facilitate learning, a skip connection copies
407 the current frame $x(t)$ and the network only outputs residuals that get added to the current frame in
408 order to predict the next frame: $\hat{x}(t+1) = x(t) + f_w([x(t), x(t-1)])$. This model jointly transforms
409 and processes pairs of frames to generate predictions, while both polar predictor (PP) and quadratic
410 predictor (QP) separate spatial processing and temporal extrapolation.

411 B Description of datasets and optimization

412 **DAVIS** To train, test and compare these models, we use the DAVIS dataset [14], which was
413 originally designed as a benchmark for video object segmentation. Image sequences in this dataset
414 contain diverse motion of scenes and objects (eg., with fixed or moving camera, and objects moving at
415 different speeds and directions), which make next frame prediction challenging. Each clip is sampled
416 at 25 frames per second, and is approximately 3 seconds long. The set is subdivided into 60 training
417 videos (4741 frames) and 30 test videos (2591 frames). We pre-processed the data, converting all
418 frames to monochrome luminance values, and scaling their range to the interval $[-1, 1]$. Frames are
419 cropped to a 256×256 central region, where most of the motion tends to occur, and then spatially
420 down-sampled to 128×128 pixels.

421 **VanHateren** We also consider a smaller dataset of natural image sequences obtained from Hans van
422 Hateren [13], as described in [3] and downloaded from <https://github.com/cadiou/twolayer>.
423 The top missing band of the images is cropped from 128 by 128 pixels to 112 by 128 pixels. The
424 dataset is standardized to zero mean and unit variance, and it is split into 292 snippets of 11 frames
425 for training and 33 snippets of 11 frames each for testing. There is no spatial downsampling or
426 whitening.

427 **Boundary handling** The computation of this prediction error is restricted to the center of the
428 image because moving content that enters from outside the video frame is inherently unpredictable.
429 Specifically, we trim a 17-pixel strip from each side, yielding frames of size 94×94 pixels.

430 **Training procedure** We assume the temporal evolution of natural signals to be sufficiently and
431 appropriately diverse for training, and do not apply any additional data augmentation procedures.
432 We train on brief temporal segments containing 11 frames, which allows for prediction of 9 frames,
433 processing these in batches of size 4. We train each model for 200 epochs on DAVIS using the
434 Adam optimizer [30] with default parameters and a learning rate of $3 \cdot 10^{-4}$. The learning rate is
435 automatically halved when the test loss plateaus. In the CNN, we use batch normalization before
436 every half-wave rectification, rescaling by the standard deviation of channel coefficients (but with no
437 additive bias terms).

438 **C Planted symmetries**

439 **D Learned filters**

440 **E Predictions examples**

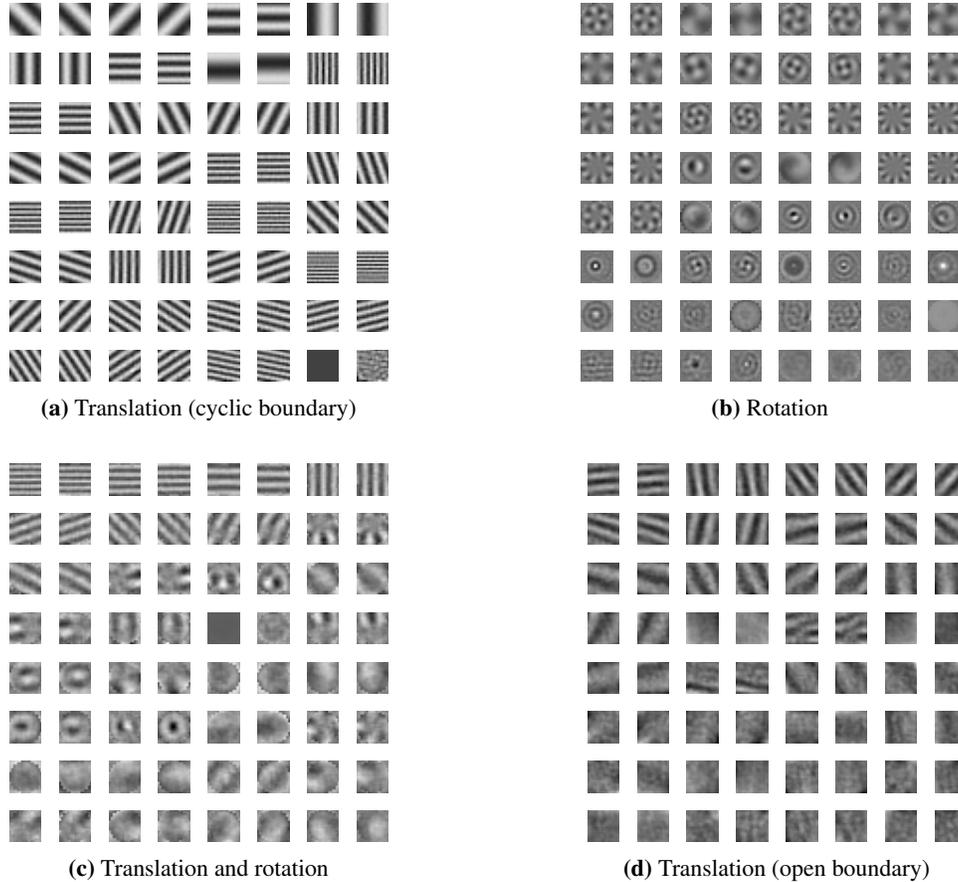


Figure 5: Filters of polar predictor networks trained to predict small synthetic sequences. We randomly select 100 image patches of size 16×16 from the DAVIS dataset and generate training data by manually transforming them - applying translations or rotations. We verify that PP recovers the known harmonic functions: Fourier modes for translation (panel **a**), and disk harmonics for rotation (panel **b**). To show that the recovery of harmonics is robust, we design two additional synthetic datasets. i) the combination of translational and rotational sequences. In this case, PP learns filters that correspond to either group, suggesting that our approach can generalize to situations with more than than one group at play (panel **c**); ii) generalized translation sequences: spatially sliding a square window on a large image (ie. new content creeps in and falls off at boundaries), instead of using cyclic boundary condition (ie. content wraps around the edges). In this case, PP learns localized Fourier-like modes (panel **d**), indicating that approximate group actions still provide meaningful training signal - although some filters are less structured. In each panel, the 32 pairs of filters are sorted by their norm. Notice that some of the filters are not structured and generally miss high frequency harmonics. This is due to the spectral properties of the datasets, which have more power at lower frequencies, and to the discretization of the transformations.

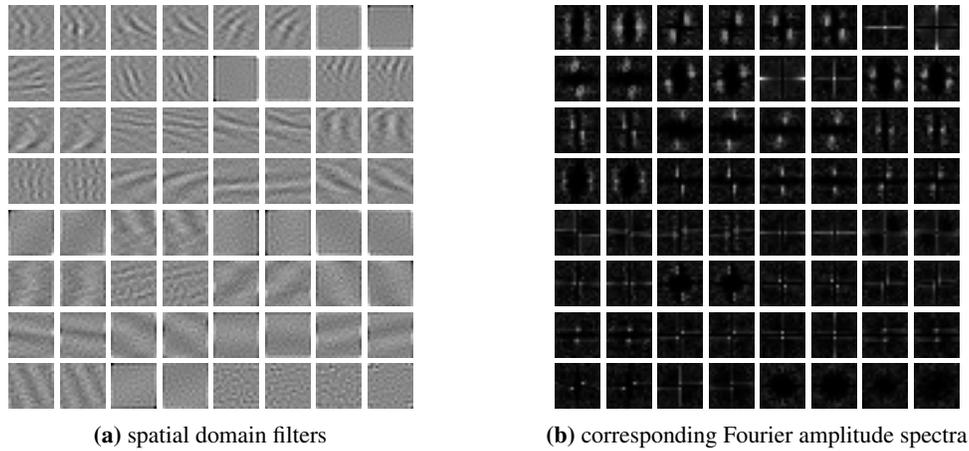


Figure 6: Filters of a polar predictor trained to predict natural videos from the DAVIS dataset. The 32 pairs of convolutional filters are sorted by their norm and their amplitude spectrum is displayed at corresponding locations on the right panel. Observe that the filters are selective for orientation and spatial frequency, tile the frequency spectrum, and form quadrature pairs.

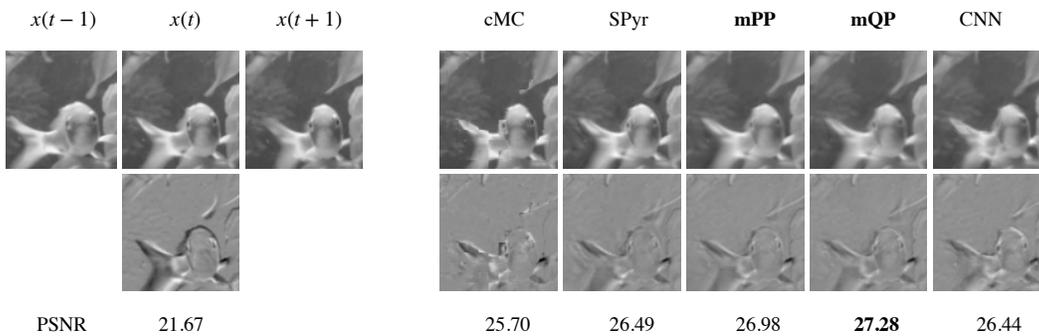


Figure 7: A typical example image sequence from the DAVIS test set. The first three frames on the top row display the unprocessed images, and last five frames show the respective prediction for each method (with their shorthand above). The bottom row displays error maps computed as the difference between the target image $x(t+1)$ and each predicted next frame on the corresponding position in the first row. Images, predictions and error maps are all shown on the same scale.