

## Supplementary Materials for **Computational multiheterodyne spectroscopy**

David Burghoff, Yang Yang, Qing Hu

Published 11 November 2016, *Sci. Adv.* **2**, e1601227 (2016)

DOI: 10.1126/sciadv.1601227

### **This PDF file includes:**

- Supplementary Materials and Methods
- fig. S1. Dual-comb spectroscopy in a chaotic regime.
- fig. S2. Cross-correction of multiple spectra.
- fig. S3. Incoherent sidebands in a weak dual-comb tooth.
- fig. S4. Multiheterodyne spectrum corrected by a two-line model.
- fig. S5. Correction of artificial dual-comb data.

## Supplementary Materials and Methods

### Nonconvexity of the error function

The error function minimized by the Kalman filter (Equation (2)) is nonconvex and consequently possesses local minima. Even if properly initialized, as a result of noise it is possible for the filter to jump into false minima, and for real-time processing it is essential that such artifacts be robustly removed. The source of these errors lies in the equivalence between the extended Kalman filter (EKF) and phase-locked-loops (PLLs): when used to track a single frequency, the EKF is a digital version of a PLL and can be thought of similarly. As a result, if two lines of the model falsely lock to two lines of the multiheterodyne signal, a poor correction will be achieved. Two types of errors are possible:

- The modeled comb's offset locks to the true offset plus an integer multiple of the true repetition rate

$$f_0^{(\text{model})} = f_0^{(\text{true})} + n\Delta f^{(\text{true})}$$

This error is fairly trivial, and is a consequence of the fact that the offset frequency of the RF comb is only defined modulo the repetition rate.

- The modeled comb's repetition rate locks to the true repetition rate times a rational number

$$\Delta f^{(\text{model})} = \frac{n}{m} \Delta f^{(\text{true})}$$

The second error is more problematic, because it affects the quality of the correction significantly. (Essentially, only two lines of the comb model lock, causing the prediction residual to increase from below 8% to over 40%.) To remedy this, an estimate for the repetition rate is pre-calculated using the coherence function  $C_\tau(t) \equiv y^*(t + \tau)y(t)$  (which contains frequency components at  $\Delta f$  in addition to its harmonics). When the filter detects that the modeled repetition rate only has strong components at two lines—signifying that the model might be falsely locked—and is also far from the true repetition rate, the filter then corrects  $\Delta f$  by multiplying by the appropriate rational number.

## Removal of the coherent artifact

Under certain conditions, a coherent artifact can be present in the extracted phase and timing signals. This artifact results in multiheterodyne spectra that *look* well-corrected, but actually have incorrect amplitudes. The reason for this is that if the process noise of the offset frequency is allowed to be large, the extracted offset frequency will often contain spurious components at harmonics of the repetition rate. These spurious components are problematic because they give rise to a frequency modulation term which cross-mixes the dual comb amplitudes. Fortunately, because the dual comb offset arises from the difference in the two combs' individual offsets (which are unrelated provided there is no strong coupling between the lasers), it should also be the case that the actual frequency component of  $f_0(t)$  at  $\Delta f$  is vanishingly small. Therefore, those components can simply be filtered out.

Dealing with this effect is not necessary when the RF comb is correctly-modeled, but is absolutely critical when the comb is mis-modeled. For example, the two-line model shown in fig. S4A is a mis-model since the filter believes that there are only two lines present, when in reality there are many more. As a result, the optimal filter would try to pull energy into the two modeled lines at the expense of the others. This would give rise to a frequency component of  $f_0(t)$  at integer multiples of  $\Delta f$ , shown in fig. S4B, and this would in turn cause the corrected comb amplitudes to be incorrect. However, note also that such effects are not present in the full model, and that at all other frequencies the two- and thirty-line models agree well. Therefore, the two-line model can be corrected heuristically by filtering out those frequency components that appear at integer multiples of the repetition rate. Such effects do not occur if the comb is *over*-modeled by including more lines than are actually present—the Kalman filter will still learn that those lines have zero amplitude—but this of course comes at the expense of increased computational complexity.

## Additional technical details

Because Kalman filters inherently process data in a serial fashion, when processed on a CPU the resources needed to correct a spectrum can be intensive. For most of the data shown here, the raw multiheterodyne data was IQ-demodulated with an intermediate local oscillator and sampled at 1.25 GS/s. Two channels were recorded, consisting of the in-phase and in-quadrature signals. (In fig. S2, two

detectors were used and so four channels were recorded.) When recorded over 100  $\mu\text{s}$ , an Intel Xeon X5680 3.33 GHz CPU can correct the data in approximately 150 s for a thirty-line model and 75 s for a two-line model.

To achieve some computational speedup, multiheterodyne data can be subdivided into sufficiently long batches, processed in parallel, and coaveraged. Ultimately, obtaining real-time performance at gigahertz sample rates will necessitate signal processing capable of handling such high data rates, and will require the development of chip-scale solutions. For example, field programmable gate arrays (FPGAs) have previously been used to continuously correct multiheterodyne data with reference signals available (11), and it is likely that a similar scheme could also be used here.

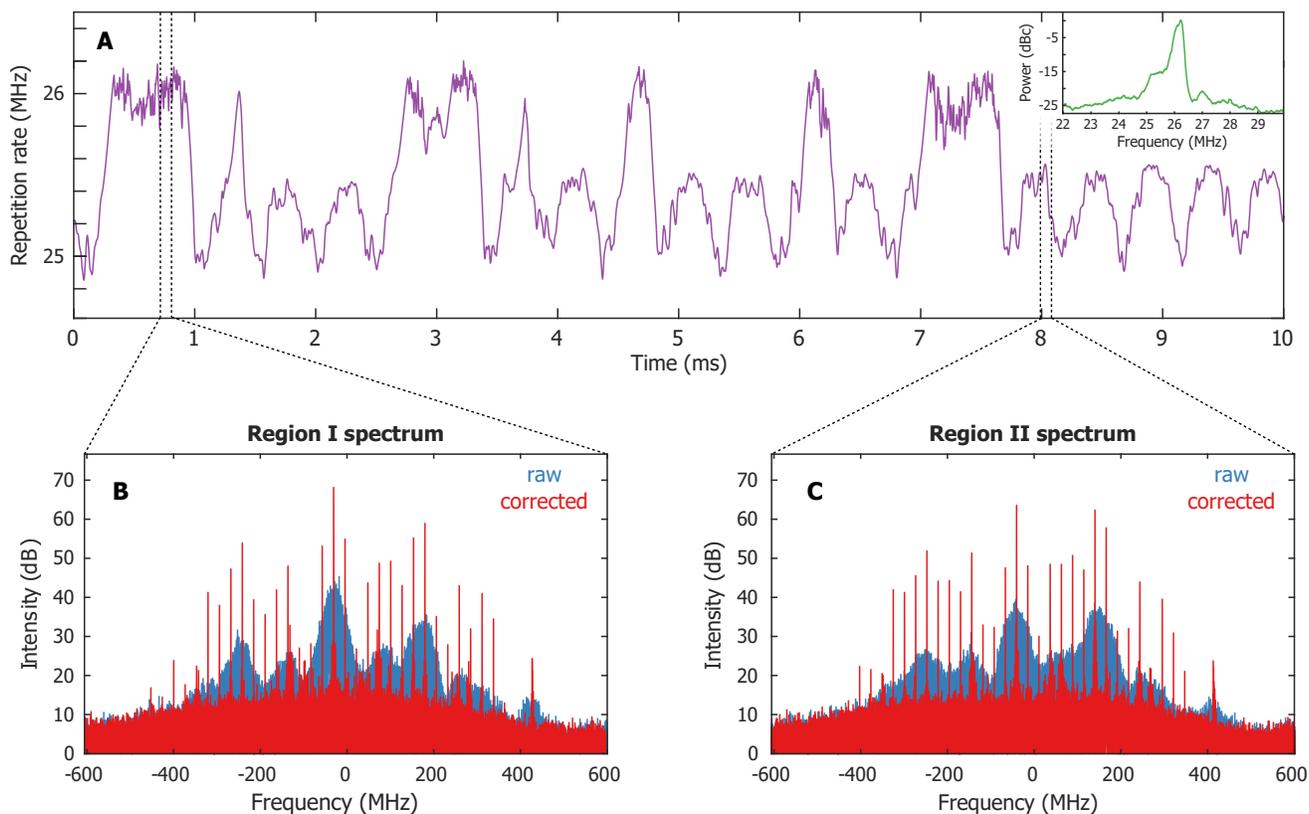
### **Correction of artificial dual comb data**

To further demonstrate the efficacy of this correction procedure and to do so in a way that does not depend on the details of quantum cascade lasers, we show in fig. S5 the results of the procedure on artificially-generated and corrupted dual comb data. Artificial dual comb data is first generated without phase or timing noise, with dual comb lines that have power levels that logarithmically span the white noise level (0 dB) up to 60 dB. The phases of each line and the order of the lines are chosen randomly, and the spectrum is plotted in green in fig. S5A. Next, phase and timing noise is introduced with a characteristic time constant of 1  $\mu\text{s}$ , and the artificial data is corrupted. The resulting multiheterodyne spectrum is plotted in blue in fig. S5A; the phase and timing errors have resulted in a completely broadened spectrum. Finally, correction is performed using *only* the information present in the broadened spectrum; the corrected spectrum (shown in red) agrees very well with the original.

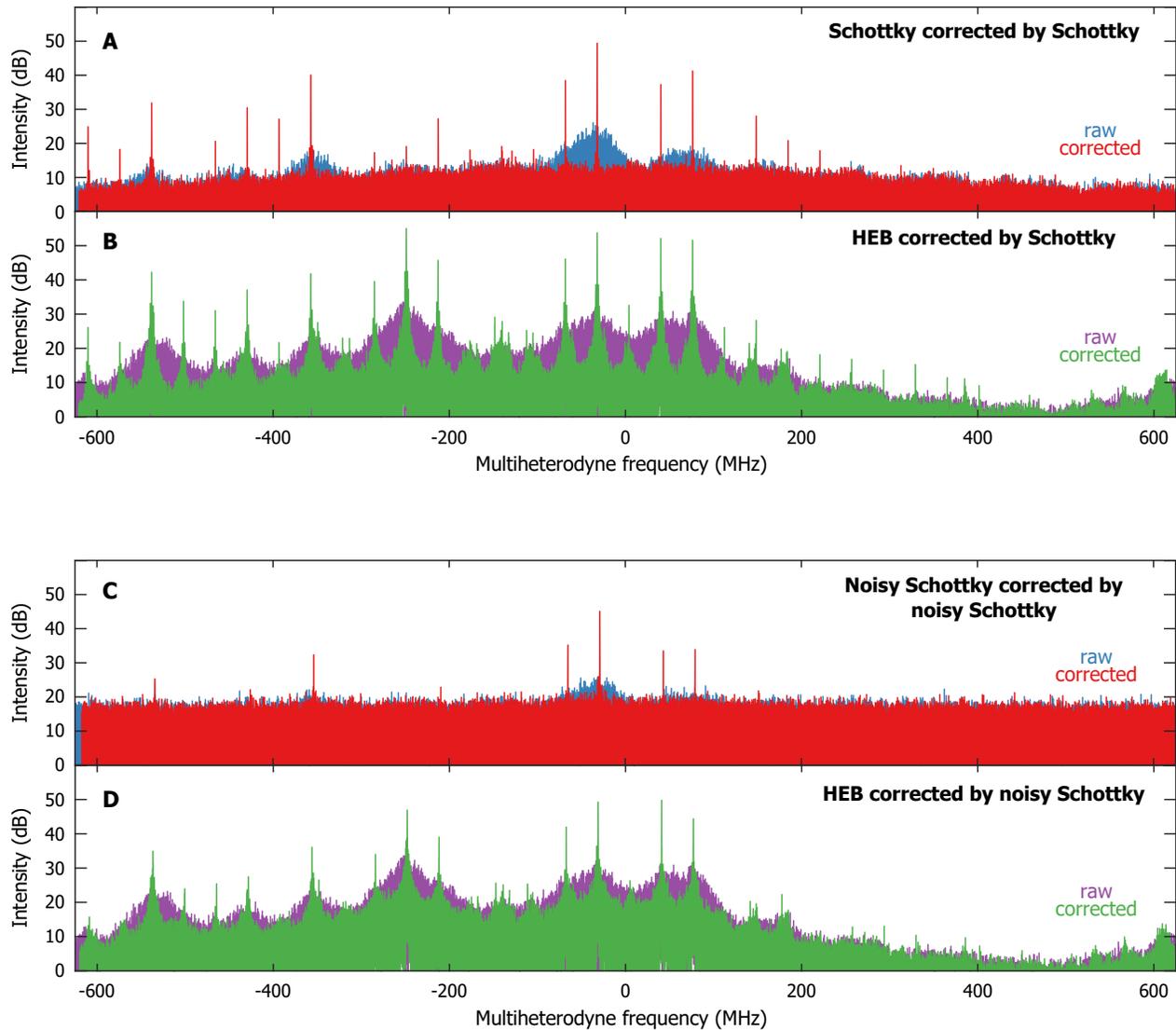
In fig. S5B, we compare the actual power level of each line (before corruption) to the estimated power level of the same line (after corruption and correction). At sufficiently high power levels the two signals agree nearly perfectly, proving that the procedure effectively preserved spectroscopic information. Figure S5C shows the same data as in fig. S5B, but in terms of fractional error instead of absolute error. The largest lines have errors under  $10^{-3}$ , with the error increasing for lines that are closer to the noise floor. Essentially, because computational correction uses some of the signal's energy to perform the

correction instead of using it to estimate the amplitude of its lines, it will always be associated with a slight reduction in the signal-to-noise ratio of the amplitude estimates. Finally, fig. S5D compares the actual phase and timing errors used to corrupt the artificial data with the phase and timing errors as estimated by the Kalman filter. Once again, good agreement is found between the two sets of signals; the filter is able to correctly track even short-term fluctuations in the two frequencies.

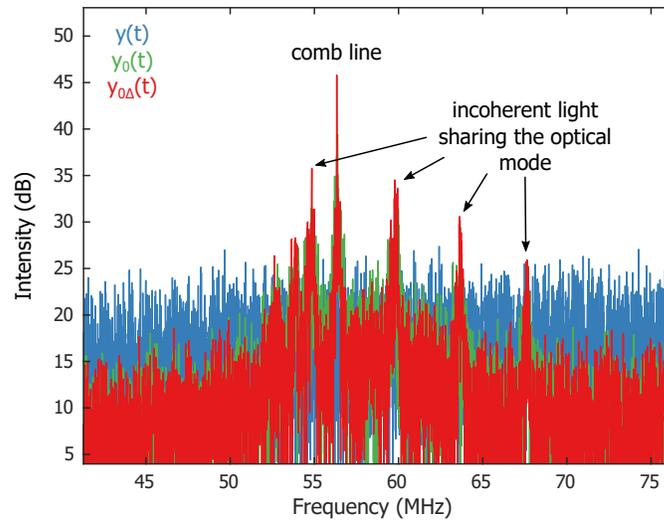
## Supplementary Figures



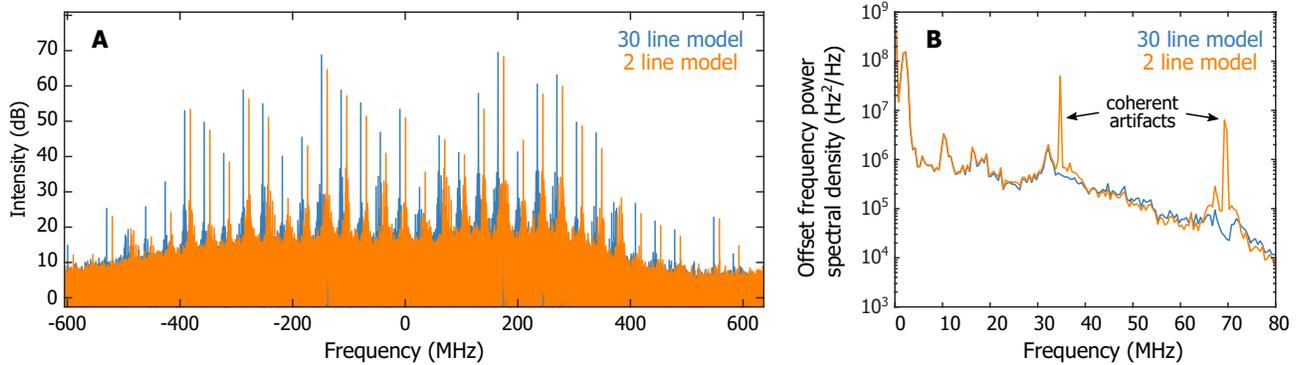
**fig. S1. Dual-comb spectroscopy in a chaotic regime.** (A) Repetition rate of a dual comb spectrum where one comb chaotically alternates between comb regimes. There is an additional 2 kHz fluctuation in the beatnote due to acoustic vibration of the cryocooler. The RF comb's repetition rate beatnote is shown in the inset. (B) and (C) Raw and corrected multiheterodyne spectra at two different comb regimes. The offset, repetition rate, and complex amplitudes all differ between regimes.



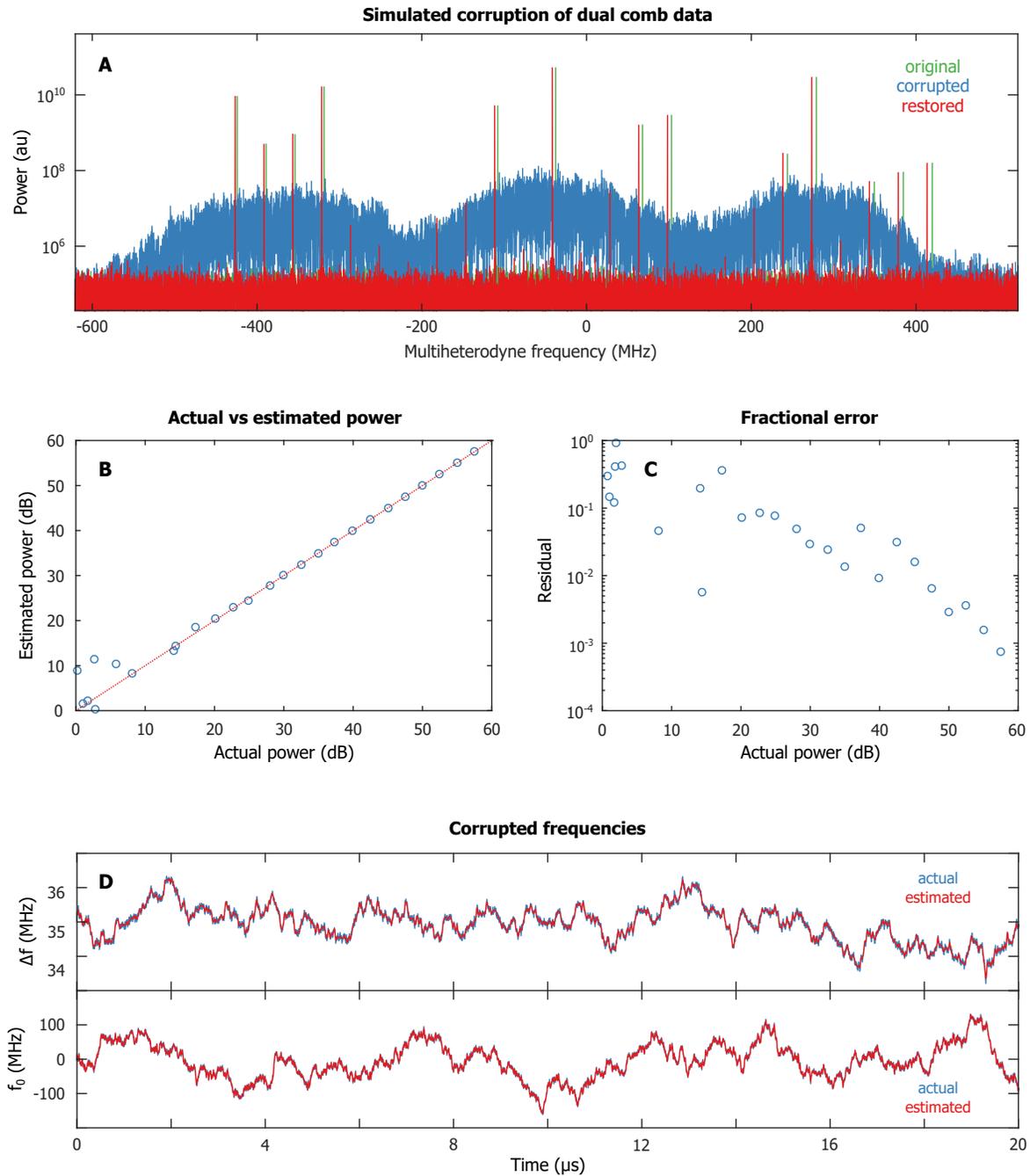
**fig. S2. Cross-correction of multiple spectra.** (A) Multiheterodyne spectrum from a Schottky mixer in which the data has been self-corrected. (B) Multiheterodyne spectrum from a hot electron bolometer (HEB) recorded at the same time, cross-corrected using error signals extracted from the Schottky mixer. (C) Same spectrum as (A), but with the raw signal corrupted by 10 dB of white noise. Even though the raw data is nearly buried in noise, correction remains possible. (D) Same spectrum as (B), but with the signal corrected by the corrupted Schottky signal. Though the correction suffers, individual teeth are still recovered.



**fig. S3. Incoherent sidebands in a weak dual-comb tooth.** Spectrum of a weak comb tooth with incoherent sidebands. Phase and timing correction cause the actual comb line to become delta function-like but only partially correct the sidebands, leaving them with substantial linewidths.



**fig. S4. Multiheterodyne spectrum corrected by a two-line model.** (A) Multiheterodyne spectrum corrected by a thirty-line model and by a two-line model, offset for clarity by 10 MHz. The two-line model reproduces the results of the thirty-line model (with significant computational savings), however the thirty-line model is more accurate and improves the SNR of most lines by several dB. (B) Power spectral density of the offset frequency  $f_0(t)$  extracted from the Kalman filter in both cases. Because the two-line model incorrectly models the system, a coherent artifact causes  $f_0(t)$  to develop components at integer multiples of the repetition rate (35 MHz here). However, such effects are absent in the thirty-line model.



**fig. S5. Correction of artificial dual-comb data.** (A) Spectrum of artificial dual comb data (green), spectrum of artificial data following phase-timing corruption (blue), and spectrum of the corrupted dual comb signal after phase-timing correction (red). (B) Comparison of the actual power of each line (before corruption) with the estimated power (after corruption and correction). The dashed line indicates perfect agreement. (C) Fractional error associated with the comparison in (B). The strongest line has a residual error under  $10^{-3}$ , and the error increases for lines that are closer to the noise floor (0 dB). (D) Comparison of the actual phase-timing errors with the phase-timing errors estimated by the Kalman filter.