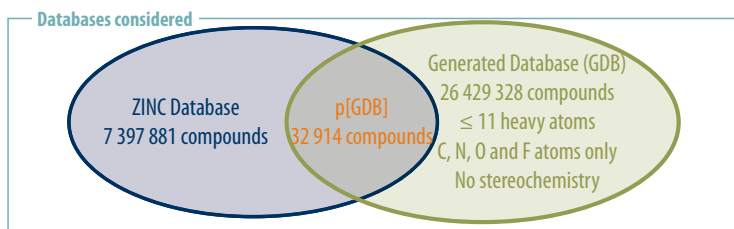
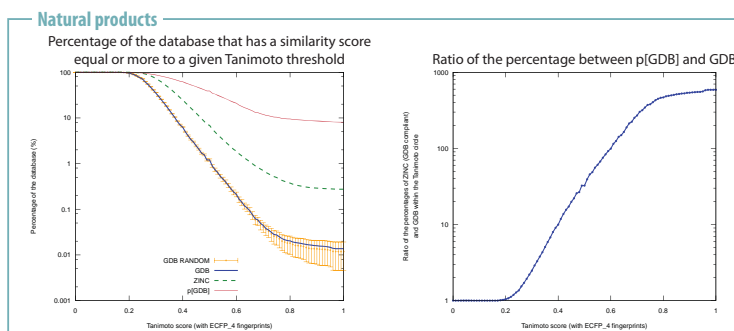
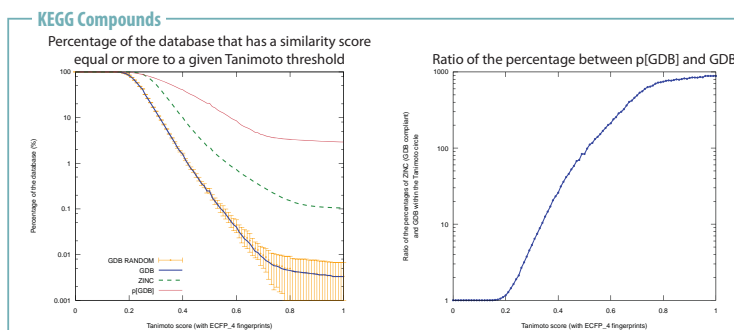


## INTRODUCTION

The size of chemical space is believed to range from  $10^{40}$  to  $10^{200}$  molecules. Why should high-throughput screening (HTS), given that it typically considers libraries of  $10^6$  molecules or less, yield any hit at all? Drug discovery nevertheless heavily relies on HTS and its ability to deliver hits and leads. We believe that screening libraries are biased toward biogenic-like molecules and propose to quantify this assumption. In what follows, chemical space, for molecules of 11 C, N, O or F atoms or less was approximated by the generated database (GDB) of Reymond *et al.*,<sup>1</sup> while screening libraries were approximated by the database of commercially available compounds (ZINC)<sup>2</sup> and its subset that meets the GDB criteria (p[GDB]).



## SCREENING LIBRARIES ARE BIASED TOWARD BIOGENIC MOLECULES

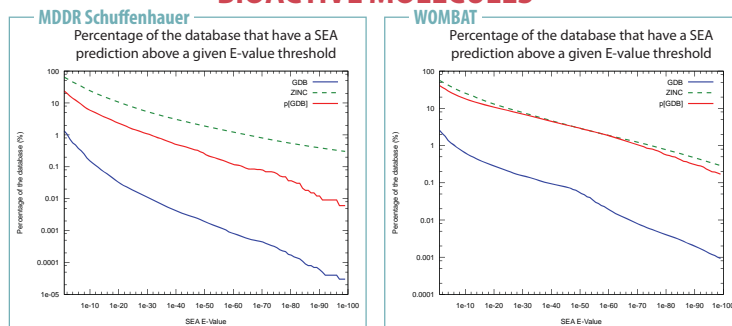


For every compound in the GDB, the p[GDB] and the ZINC database, the maximum similarity to KEGG compounds<sup>3</sup> and natural products<sup>4</sup> was calculated. The left figures show the percentage of the database that had a similarity score equal or better than a given Tanimoto threshold. When KEGG compounds were the reference database, only 5% of the compounds in the GDB database have a similarity equal or higher than 0.34 compared to 57% for the p[GDB] database. At this Tanimoto threshold, the p[GDB] database was hence 10 times as similar to metabolic space as the GDB database. This ratio steadily increased with the Tanimoto threshold to reach almost 3 orders of magnitude between p[GDB] and GDB (right figures). The purchasable compounds were biased toward biogenic molecules: their distribution would be comparable to that of the GDB if they were charting random parts of chemical space (orange lines on the right figures).

## REFERENCES

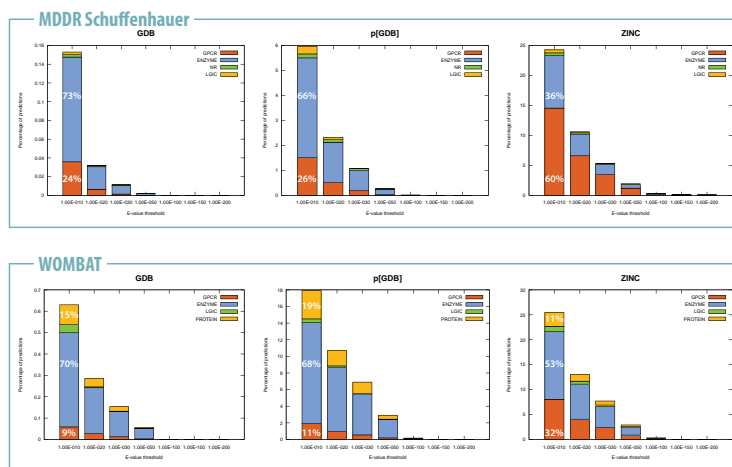
1. T. Fink and J. L. Reymond, *J. Chem. Inf. Model.* 47 (2), 342 (2007).
2. ZINC is available from <http://zinc.docking.org>.
3. KEGG is available from <http://www.genome.jp/kegg>.
4. The Dictionary of Natural Products is available from <http://dnp.chemnetbase.com>.
5. M. J. Keiser, B. L. Roth, B. N. Armbruster *et al.*, *Nat. Biotech.* 25 (2), 197 (2007).
6. The MDL Drug Data Report database is available from <http://mdl.i.com>
7. M. Olah *et al.*, in *Cheminformatics in Drug Discovery*, edited by T.I. Oprea (Wiley-VCH, New York, 2004), pp. 223.

## SCREENING LIBRARIES ARE BIASED TOWARD BIOACTIVE MOLECULES



SEA<sup>5</sup> scores were calculated for every compound in the GDB, the p[GDB] and the ZINC database against the ligand sets in the MDDR Schuffenhauer<sup>6</sup> and the WOMBAT databases.<sup>7</sup> 6% of the compounds in the p[GDB] had an E-value of  $10^{-10}$  or better to one of the sets in the MDDR databases; in contrast only 0.15% of the compounds in the GDB had a prediction. The ratio of compounds that had a prediction at a given E-value threshold easily reached 2 orders of magnitude between the p[GDB] and GDB databases.

## SCREENING LIBRARIES ARE BIASED TOWARD DRUGGABLE TARGETS



For the GDB database, 73% and 70% of the targets that were predicted using SEA (with an E-value  $\leq 10^{-10}$ ) were enzymes while only 24% and 9% were GPCRs, for the MDDR and the WOMBAT databases, respectively. A similar distribution of target types is observed for the predictions of the p[GDB] compounds. In contrast, 60% and 32% of the targets predicted for the ZINC compounds were GPCRs. Although the p[GDB] database is more similar to biogenic and bioactive molecules than the GDB database it does not bear the same bias toward GPCR like compounds as the ZINC compounds.

## WHAT'S MISSING FROM p[GDB] AND GDB?

