Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

# Reliability of Causality Assessment for Drug, Herbal and Dietary Supplement Hepatoxicity in the Drug-Induced Liver Injury Network (DILIN)

**Paul H. Hayashi, MD, MPH**[1],[*], **Huiman X. Barnhart, PhD, MS**[2],[**], **Robert J. Fontana, MD**[3],[*], **Naga Chalasani, MD**[4],[*], **Timothy J. Davern, MD**[5],[*], **Jayant A. Talwalkar, MD**[6],[*], **K. Rajender Reddy, MD**[7],[*], **Andrew A. Stolz, MD**[8],[*], **Jay H. Hoofnagle, MD**[9],[*], and **Don C. Rockey, MD**[10],[*]

[1]Division of Gastroenterology and Hepatology, University of North Carolina, Chapel Hill, NC

[2]Duke Clinical Research Institute & Department of Biostatistics, Duke University, Durham, NC

[3]Division of Gastroenterology, University of Michigan, Ann Arbor, MI

[4]Division of Gastroenterology and Hepatology, Indiana University School of Medicine, Indianapolis, IN

[5]Division of Gastroenterology, California Pacific Medical Center, San Francisco, CA

[6]Division of Gastroenterology and Hepatology, Mayo Clinic, Rochester, MN

[7]Division of Gastroenterology, University of Pennsylvania, Philadelphia, PA

[8]Division of Gastrointestinal and Liver Diseases, University of Southern California, Los Angeles, CA

[9]Liver Disease Research Branch, Division of Digestive Diseases and Nutrition, National Institute of Diabetes and Digestive and Kidney Diseases (NIDDK), NIH, Bethesda, MD

[10]Division of Gastroenterology, University of Texas Southwestern, Dallas, TX

## Abstract

**Background**—Due to the lack of objective tests to diagnose drug induced liver injury (DILI), causality assessment is a matter of debate. Expert opinion is often used in research and industry but its test-retest reliability is unknown.

[*]Author Contributions:
The author contributed to the study concept and design, acquisition of data, analysis and interpretation of data, drafting of the manuscript, and critical revision of the manuscript for important intellectual content.

[**]The author contributed to the study concept and design, acquisition of data, analysis and interpretation of data, critical revision of the manuscript for important intellectual content, and biostatistical analysis.

**Corresponding author:** Paul H. Hayashi, MD, MPH, UNC Liver Program, Division of Gastroenterology & Hepatology, CB# 7584 Burnett-Womack Bldg., Room 8011, Chapel Hill, NC 27599-7584, Phone 919-966-2516, Fax 919-966-1700, paul_hayashi@med.unc.edu.

Disclosures:
PH Hayashi: No relevant disclosures or conflicts of interest; HX Barnhart: No relevant disclosures or conflicts of interest; RJ Fontana: Consultant GSK, Tibotec, Vertex, Gilead; N Chalasani: Consultant Merck, Abbvie, Aegerion, Salix, BMS, Lilly; TJ Davern: No relevant disclosures or conflicts of interest; JA Talwalkar: No relevant disclosures or conflicts of interest; KR Reddy: Advisory Board BMS, Abbvie, Gilead, Merck Genetech-Roche, Vertes, Jannsen; AS Stolz: No relevant disclosures or conflicts of interest; JH Hoofnagle: No relevant disclosures or conflicts of interest; DC Rockey: No relevant disclosures or conflicts of interest

**Aims—**To determine the test-retest reliability of the expert opinion process used by the Drug-Induced Liver Injury Network (DILIN)

**Methods—**Three DILIN hepatologists adjudicate suspected hepatotoxicity cases to 1 of 5 categories representing levels of likelihood of DILI. Adjudication is based on retrospective assessment of gathered case data that includes prospective follow-up information. One hundred randomly selected DILIN cases were re-assessed using the same processes for initial assessment but by 3 different reviewers in 92% of cases.

**Results—**The median time between assessments was 938 days (range: 140–2352). Thirty-one cases involved >1 agent. Weighted kappa statistics for overall case and individual agent category agreement were 0.60 (95% CI: 0.50–0.71) and 0.60 (0.52–0.68), respectively. Overall case adjudications were within one category of each other 93% of the time, while 5% differed by 2 categories and 2% differed by 3 categories. Fourteen-percent crossed the 50% threshold of likelihood due to competing diagnoses or atypical timing between drug exposure and injury.

**Conclusions—**The DILIN expert opinion causality assessment method has moderate inter-observer reliability but very good agreement within 1 category. A small but important proportion of cases could not be reliably diagnosed as 50% likely to be DILI.

## INTRODUCTION

In 2004, the National Institute of Diabetes and Digestive and Kidney Diseases (NIDDK) established the Drug-Induced Liver Injury Network (DILIN) as a multi-center study with aims of improving the understanding of the causes, outcomes and molecular mechanisms of hepatotoxicity due to medications or herbal and dietary supplements (HDS)(1). The Network of 8 centers has enrolled >1200 patients into the Prospective Study. Cases meeting laboratory enrollment criteria are enrolled at the discretion of the investigators at each site and based on their clinical suspicion that such abnormalities are at least possibly due to DILI. Each case has a clinical narrative and a formatted file of demographics, medications, radiography data, histology (when available), and laboratory values in flow sheet format. Serum, plasma and DNA are also collected.

Because there are no objective diagnostic tests, DILI remains a diagnosis of exclusion that requires adequate clinical, laboratory and imaging data. Scoring algorithms for diagnosis are available (2–4) but retest reliability studies are quite limited. Roussel Uclaf Causality Assessment Method (RUCAM) that was developed under the auspices of the Council for International Organizations of Medical Sciences (CIOMS) is the most widely accepted and validated instrument, yet only one study has examined its retest reliability (5). For research purposes, DILIN uses a standardized procedure for expert opinion consensus to adjudicate the likelihood of DILI that is based in part on RUCAM and yields similar categories of DILI likelihood.(6) However, reliability of the DILIN expert opinion process has not been critically assessed. Poor reliability would undermine any future mechanistic studies using serum, plasma or DNA from the DILIN subjects and undermine the use of DILIN cases for

development of other more accessible diagnostic algorithms for the clinician. Therefore, the aim of this analysis was to assess the inter-rater, test-retest reliability of the DILIN consensus opinion process.

## METHODS

### DILIN Prospective Cohort

The DILIN study has been previously described in detail.(6) Patients suspected of having liver injury due to medications or HDS products were enrolled within 24 weeks of injury onset and then followed prospectively for 6 to 24 months depending on the pace and completeness of DILI resolution. Because the enrollment window was 24 weeks, cases enrolled at varying time points in their DILI event. Enrollment criteria were (1) serum alanine aminotransferase (ALT) or aspartate aminotransferase (AST) levels > 5 times the upper limit of the normal (ULN) (or pretreatment baseline if abnormal) on 2 consecutive occasions, or (2) alkaline phosphatase (AP) levels > twice the ULN (or pretreatment baseline if abnormal) on 2 consecutive occasions, or (3) total serum bilirubin > 2.5 mg/dL, or international normalized ratio (INR) > 1.5 with any elevation in serum ALT, AST or AP. Cases meeting laboratory enrollment criteria are enrolled at the discretion of the investigators at each sight and based on their clinical opinion. Case enrollment is not restricted by time of onset because some medications are well known to have very long latencies of even years (e.g. nitrofurantoin). We did not restrict enrollment to those with full follow-up data showing resolution of injury (i.e. dechallenge) because such restriction would hinder the prospective nature of DILIN. We wanted to capture this dechallenge data while under study protocol. Injury pattern was categorized as cholestatic (R<2), mixed (R 2–5) or hepatocellular (R>5) where $R = (ALT/ULN) \div (AP/ULN)$. Severity level was based on INR, bilirubin, signs of liver failure, need for hospitalization and fatal or transplant outcome as previously described.(6) Exclusion criteria included acetaminophen hepatotoxicity, prior liver or bone marrow transplant, alcohol related liver disease, autoimmune hepatitis or genetic liver disease. Patients with compensated chronic hepatitis B, C or with nonalcoholic fatty liver disease were eligible and enrolled at the discretion of the site investigators. For patients with such background liver disease, reviewers used baseline liver enzyme levels, viral serologies and viral nucleic acid tests to judge the presence of DILI versus exacerbation of underlying liver disease.

At baseline visit, a detailed history was obtained, and clinical, laboratory, and imaging results extracted from records. As reported previously, a DILIN protocol battery of tests to exclude other causes of liver injury were obtained at enrollment if not extractable from chart records (6). Serum, plasma, urine and DNA specimens were sent to a central repository for future studies. Subjects were followed for at least 6 months and those with persistent liver abnormalities or signs of chronic liver injury were followed through 24 months.

Ninety (90%) of the cases had complete documentation of data for all 21 parameters adapted from Agarwal et al. as essential data for DILI cases.(7) (Appendix figure) Nine (9%) were missing complete data for one parameter each (4 cases -- incomplete viral serologies, 4 cases – no documented hepatic imaging, 1 case -- incomplete 'washout' of liver biochemistries); one case (1%) was missing complete data for 2 parameters, viral serologies and autoimmune

markers. This last case was considered unlikely to be DILI in large part due to lack of these data. Overall, the 90 cases with complete data were considered more likely to be DILI than the 10 cases with incomplete data (median DILIN scores of 2 versus 3, respectively, p = 0.04).

## DILIN Expert Opinion Process

DILIN causality assessment has been described in detail (6). The process is the same for medication and HDS hepatotoxicity, since there is little data to suggest a different causality process is necessary or valid. Each case was adjudicated independently by 3 hepatologists including the site investigator who enrolled the case. Assessment is based on retrospective review of the case history, laboratory data and prospective follow-up study visits. Each reviewer assigns a causality score corresponding to percentages of DILI likelihood in which 1= definite (>95% likelihood), 2 = very likely (75–95%), 3 = probable (50–74%), 4 = possible (25–49%), and 5 = unlikely (<25%). Disagreements were identified after the 3 reviewers submitted their scores. Consensus scores were achieved by electronic mail or conference call discussions. Those cases in which agreement could not be reached by the three reviewers were then voted upon by one member from each DILIN site during monthly conference calls. Final score was assigned by majority vote. For cases involving >1 agent, an overall case score and separate individual scores were determined. For example, the overall case score might be 1 (definite DILI) with one agent scoring a 2 (very likely causal) and the other scoring a 4 (only possibly).

From inception in 2004 to April 16, 2009, adjudication was based upon results obtained shortly after enrollment thus simulating the clinician's task to assess at time of presentation. However, when enrollment occurred within in days of onset, data on resolution of injury was sparse. After April 16, 2009, the protocol was changed so that adjudication was done 6 months after enrollment so that follow-up data could be included into the assessment.

## Reliability Cohort and Reassessment

The DILIN Data Coordinating Center chose 100 cases by computer driven random assignment from the Prospective registry. Two cases involving the interval development of new diagnostic information regarding hepatitis E testing were included but the HEV data was specifically excluded for reassessment.(8) Chosen cases were stratified 1:1 across April 16, 2009. Group A included 49 cases enrolled before April 16, 2009, and Group B, 51 cases enrolled afterwards. Group A cases did not have 6-month data for the initial assessment but these data were available for the reassessment. Group B cases had 6-month data for both the initial and reassessments. We stratified across these two periods to examine whether reliability is influenced by using follow-up data. For reassessments, 92 cases had 3 new reviewers. Due to an administrative error, 8 cases had one previous reviewer and 2 new reviewers. No cases were reassessed by the site investigator who enrolled the case. The rationale for excluding the enrolling site investigator from reassessments was to minimize recall bias as the enrolling investigator is often the hepatologist who continues to care for the patient. At least 4 months had to elapse before a case could be selected for reassessment. The process of reassessment was otherwise the same.

## Analysis

Standard descriptive statistics were used to characterize this cohort of 100 patients and compared to the remaining 983 cases not selected. The original and re-adjudication scores were compared for the 100 cases using weighted kappa statistics due to ordinal nature of the categories. Because some cases involved >1 drug or HDS product, the reliability at the individual agent level was also examined. Weighted kappa statistics were also determined for Group A and Group B cases separately. Because the two categories of highest likelihood (>95% likelihood and 75–95%), were similar in clinical and research relevance, reliability was assessed collapsing these two categories into one thus creating quartiles of percent likelihood. Score differences were categorized as 1, 2 or >2 scores apart and direction of changes were tallied. Cases with scores crossing the 50% likelihood line (1, 2, or 3 vs. 4 or 5) were re-examined in detail to assess factors that might have led to such uncertainty in DILI diagnosis.

# RESULTS

## Subjects

Among 1083 DILIN cases adjudicated by June 1, 2011, 49 assessed prior to the use of 6-month follow-up data (Group A) and 51 assessed with the use of 6-month data (Group B) were randomly chosen. These cases were similar to those not chosen across a variety of clinical and demographic variables (Table 1). Seventeen of the 100 had initial agreement without need for any discussion whatsoever. This rate is similar to the 20.1% for the total cohort. Of the 100 cases, 69 involved only one medication or HDS product and 31 involved multiple agents. A total of 138 different agents were implicated. The original causality scores were similar between Groups A and B (Table 2).

## Reliability of the DILIN Expert Opinion Causality Assessment Process

The median time between initial and re-assessment was 938 days (range: 140–2352). Cross tabulation of scores between original assessment and re-assessment are shown in Table 2a. Weighted kappa statistic for score agreement for overall DILI diagnosis was 0.60, 95% confidence interval (CI) 0.50–0.71. Kappa for the 69 single agent cases and 138 individual agent scores were similar [0.59 (0.46–0.71) and 0.60 (0.52–0.68)]. Score agreement tended to be better for Group B compared to Group A for overall case scores [0.67 (0.38–0.81) vs. 0.53 (0.38–0.68)], single agent cases [0.69 (0.53–0.85) vs. 0.49 (0.31–0.68)] and individual agent scores [0.66 (0.56–0.76) vs. 0.49 (0.34–0.63)], although the differences did not reach statistical significance. Collapsing score categories 1 and 2 into one category did not change the kappa score for overall case score (0.60 [0.48–0.72]) but the kappa for Group B increased to 0.73 (0.59–0.87) while the kappa for Group A fell to 0.44 (0.27–0.61) (Table 2b). Excluding the 8 cases that had one repeat reviewer on reassessment did not change the kappa scores significantly (data not shown).

## Magnitude of disagreement

The magnitude of disagreement on reassessment was small. 93% of overall scores were the same or differed by only 1 point (92% for Group A; 94% for Group B). (Figure 1a). When

scores 1 and 2 were combined into one category, 95 of 100 were within one score (92% for Group A; 98% for Group B). (Figure 1b)

### Direction of disagreement

Overall causality reassessment scores had a lower likelihood of DILI compared to the initial evaluation scores, with a reduction in the average initial score – average reassessment score of −0.22 (range 3 to −2). Of the 45 cases where scores differed, 15 were considered more likely to be DILI, while 30 were considered less likely to be DILI on reassessment. Direction of changes was similar between groups (Group A: 7 more likely, 17 less likely; Group B: 8 more likely, 13 less likely; p = ns).

Because there were more cases originally scoring 1–2 (62%) than 4–5 (19%) reflecting the careful consideration of other potential causes for liver injury before enrollment, the DILIN cohort was prone to a ceiling effect. Therefore we looked at cases originally scoring in the middle as well as proportionate increases versus decreases in the reassessment scoring of cases originally scoring a 2 or 4 respectively. Cases initially scored as a 3 still tended to score less likely to be a DILI event (42% vs. 32%), and cases originally scoring 2, were reassessed as less likely to be DILI proportionately more often than cases scoring 4 were reassessed as more likely (33% versus 20%). (Table 3)

### Cases crossing the 50% likelihood on reassessment

There were 14 (14%) cases with scores that crossed the 50% likelihood threshold on re-assessment (scores 1–3 vs. 4–5). Eight cases were from Group A and 6 from Group B (p = ns). In these 14 cases, 12 scores crossed below the 50% threshold while only 2 crossed the line toward more likely DILI (Table 4). Nine of 14 (64%) differed by one point only (i.e. between 3 and 4), but 5 differed by >1 point. Eleven involved only one drug or HDS product. No implicated agents appeared more than once in these 14 cases. Although not statistically significant (p = 0.51), the reassessment score differed from the site investigator's initial assessment score more often (12 of 14) than it did for the other two initial reviewers (9 of 14 for both).

Review of these 14 cases and recorded comments by reviewers suggest two major reasons for diagnostic uncertainty. Four cases had uncertain or inconsistent timing between agent exposure and liver injury, 7 had competing diagnoses and 1 had both reasons (Table 5).

## DISCUSSION

The DILIN expert opinion process for causality assessment has moderate test-retest reliability on a 5-point scale of likelihood, but agreement within one category of likelihood was very good at 93% of cases. Inclusion of 6-month follow-up data tended to improve concordance, suggesting diagnostic reliability probably improves when longer follow-up is available. These data are critical in establishing the DILIN registry as a source of tissue and serum for mechanistic studies. These data also help establish the registry as a reliable source of cases for the development of diagnostic instruments that are clinically accessible. While consensus expert opinion may be a reliable diagnostic method, it is cumbersome, inaccessible to the clinician and used only for research purposes. Development of a

computerized diagnostic tool for clinicians will need large registries of reliably diagnosed cases.(9)

The reliability of the DILIN consensus process is better than that reported for individual assessments without consensus. Studies of reliability between individuals diagnosing adverse drug reactions of all types yield kappa statistics as low as 0.05 to 0.2. The consensus process used by DILIN substantially elevates the inter-rater reliability (kappa 0.60). The consensus of 3 reviewers per case attenuates individual biases and variations in experience. Moreover, the collective experience and expertise of the larger causality committee (15–25 hepatologists per call) is brought to bear on cases where the 3 reviewers cannot agree. Thus, the consensus process elevates the test-retest reliability to a level comparable to that reported for histologic liver diagnoses that also rely on consensus expert opinion. Chronic viral hepatitis biopsies reviewed by 4 expert liver histopathologists produced kappa statistics for disease activity and fibrosis of 0.43 and 0.59 respectively.(10) Similar reliability results were observed for the histologic diagnosis of NASH (kappa of 0.61).(11) Therefore, if DILI diagnosis by expert opinion is to be considered as reliable as interpretation of liver biopsies, our data suggest that a rigorous consensus process must be incorporated.

Reliability measured by weighted kappa and percent agreement within one score were consistently better in Group B for overall case and individual agent scores, although the differences did not reach statistical significance probably due to small sample size. Group A's lack of requirement of 6-month follow-up data is a significant handicap since most medical diagnoses become more reliable over time. In addition, operational variability probably decreased over time as the DILIN reviewers became more familiar with the scoring scale and consensus process. For these reasons, Group B more accurately reflects the DILIN's current reliability. When categories 1 and 2 were combined, the kappa for Group B improved to 0.73 with 98% of case reassessments being within one score of each other.

Cases originally scoring in the middle of the scale at 3 were most likely to have different scores on reassessment (Table 4) probably due to less certainty of DILI diagnosis as well as being midway in a 5-point scale. Reassessment scoring tended to decrease as opposed to increase in likelihood of DILI (Table 4). Exclusion of the enrolling site investigator from reassessment may have contributed to this finding, because the site investigator was often the hepatologist who also provided clinical care to the patient. Such firsthand knowledge of the case may provide more accurate causality assessment and thereby provide the site investigator a stronger position from which to advocate for a particular score during discussions. Nevertheless, the consensus process is robust enough to still produce reasonable kappa values and very close agreement within one score without firsthand knowledge of the cases on reassessment.

As with any diagnostic tool, DILIN expert opinion is prone to interval discoveries. Recently a small percentage of DILI cases were discovered to have evidence of acute hepatitis E virus (HEV) infection. (8, 12). Two patients included in this study had data on hepatitis E infection discovered between the two assessments. For studying retest reliability of the DILIN *processes* only, we expunged this HEV data from the documents reviewed by the second set of reviewers. However, accuracy mandates that such interval discoveries be

continually incorporated and cases reassessed as were done for the DILIN HEV cases.(8) This reassessment process and updating of scores is built into the DILIN protocol. Concerns over such changes should not halt efforts to accrue cases and model diagnostic instruments using the DILIN registry for training and validation purposes. Ideally, any new diagnostic models for clinicians should be malleable enough to incorporate new discoveries as they become available.

DILIN cases vary in complexity particularly in regards to whether or not competing causes of liver injury are identified. The 14 (14%) cases that could not be reliably diagnosed as DILI with at least 50% certainty (i.e. cases crossing between 3 and 4 on reassessment) were some of the more complicated cases that had equivocal presenting diagnostic or longitudinal data. Because these complex cases are a part of clinical practice, a detailed examination of them may guide the building of an accessible diagnostic instrument for clinicians. None of these 14 cases involved isoniazid (INH) or amoxicillin-clavulanate though these drugs were the two most frequently implicated in the DILIN registry accounting for >15% of all cases. (13) The well-established signature patterns of injury for these agents (14–16) probably make DILI easier to adjudicate on one side or the other of the 50% threshold. Giving more weight to such signature presentations may improve the reliability of future diagnostic algorithms.

An alternative cause of liver injury was identified in 10 of the 14 cases and many of these alternative diagnoses lack objective, confirmatory tests (e.g. ischemic hepatitis, alcohol-related liver injury). Alternative diagnoses are well known to complicate attribution of a liver injury to a specific medication or herbal product.(17–18) Incorporation of diagnostic criteria for competing disorders (e.g. International Autoimmune Hepatitis Group diagnostic criteria (19) in expert opinion and future causality instruments may improve reliability in such cases. Expert opinion struggled with missing of precise data on timing of liver enzyme abnormalities and agent exposure in 5 cases. Perhaps, clear-cut timing of agent start and stop and enzyme elevation should be a minimum requirement for assessment akin to minimum requirements suggested for DILI case reporting.(20) Rare or unknown hepatotoxicity may have contributed to uncertainty in 2 cases (ranitidine and an experimental agent), so it would be useful to have standardized scoring of published data that is more precise than what is found in RUCAM. The LiverTox on-line textbook developed by the NIDDK, National Library of Medicine, and DILIN contains an extensive listing of publications related to DILI attributed to several hundred agents (http://livertox.nih.gov/index.html). This site may prove useful in standardizing assessment of published data on a particular agent.

Progress in the prevention, early detection and treatment of DILI will require well-characterized and prospectively followed cases of injury attributed to specific agents. Only with large registries of reliably diagnosed DILI will progress be made in determining the molecular mechanisms of DILI. This analysis suggests that the expert opinion causality assessment process used in DILIN will provide a cohort in whom the majority of cases are reliably diagnosed as DILI or not. Specifically, the updated and now well-practiced causality assessment process yields moderate diagnostic reliability based on kappa statistics and excellent agreement within one score.

While such internal reliability is encouraging, external reliability amongst non-DILIN experts using the same consensus process would be worth examining. The DILIN experience is limited to the U.S., and therefore applicability to other countries is unclear. The DILIN expert opinion process also lacks quantitative scores of individual characteristics that may always hinder retest reliability Importantly, our study does not address the daunting problem of validity since there is no gold standard for the diagnosis of DILI. Some have suggested using the RUCAM followed by expert opinion in a two-step process to enhance reliability and validity.(21)

Finally, a minority of cases straddles the 50% likelihood line and eludes a reliable diagnosis of DILI versus not DILI usually because of lack of typical presentation, imprecise information on timing or the presence of competing causes of liver injury. Cases scoring in this middle range in general, but particularly when concerns over timing or competing diagnoses are raised, will need to be reviewed carefully if used for mechanistic studies and deserve special attention if a more automated and widely accessible causality assessment instrument is to be developed.

## Acknowledgment

## Abbreviations

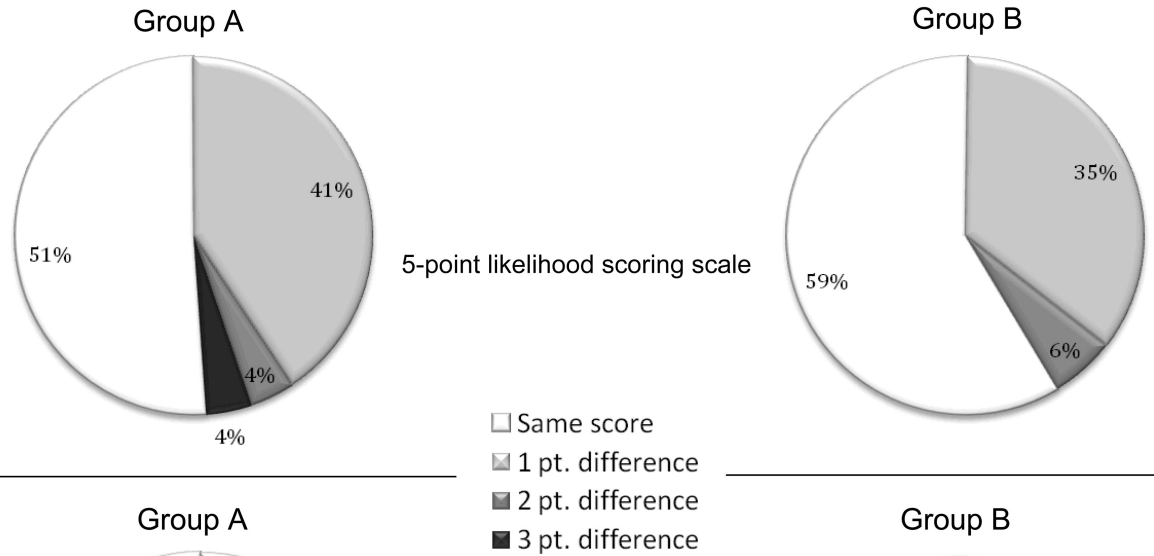| | |
|---|---|
| **ALT** | alanine aminotransferase |
| **AMA** | Anti-mitochondrial antibody |
| **ANA** | Anti-nuclear antibody |
| **AP** | Alkaline phosphatase |
| **APAP** | Acetaminophen |
| **ASMA** | anti-smooth muscle antibody |
| **AST** | aspartate aminotransferase |
| **DILI** | Drug-induced Liver Injury |
| **DILIN** | Drug-induced Liver Injury Network |
| **HDS** | Herbal and Dietary Supplements |
| **LFTs** | Liver Function Tests (AST, ALT, AP, bilirubin) |

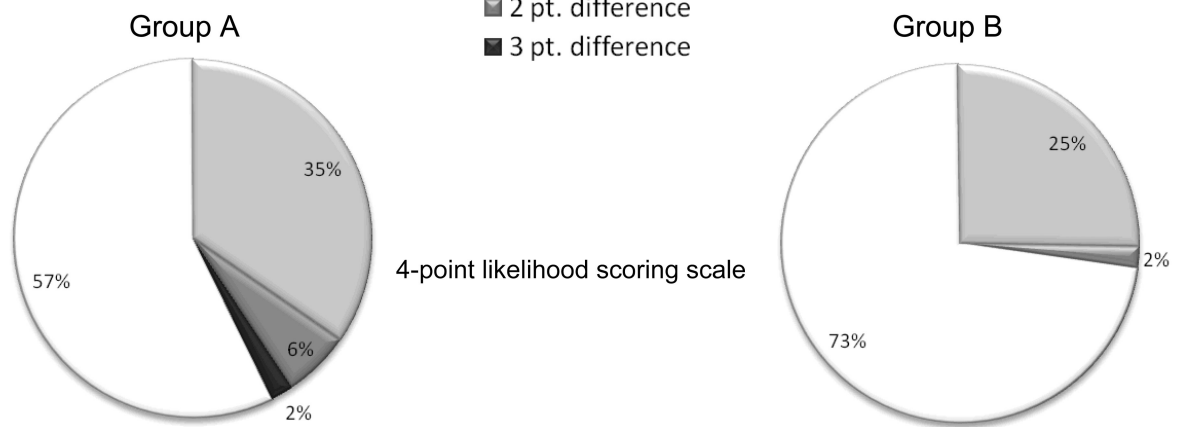| NIDDK | National Institute of Diabetes and Digestive and Kidney Diseases |
| NIH | National Institutes of Health |
| RUCAM | Roussel Uclaf Causality Assessment Method |
| ULN | upper limit of normal |

## REFERENCES

1. Hoofnagle JH. Drug-induced liver injury network (DILIN). Hepatology. 2004 Oct.40(4):773. [PubMed: 15382161]

2. Danan G, Benichou C. Causality assessment of adverse reactions to drugs--I. A novel method based on the conclusions of international consensus meetings: application to drug-induced liver injuries. J Clin Epidemiol. 1993 Nov; 46(11):1323–1330. [PubMed: 8229110]

3. Benichou C, Danan G, Flahault A. Causality assessment of adverse reactions to drugs--II. An original model for validation of drug causality assessment methods: case reports with positive rechallenge. J Clin Epidemiol. 1993 Nov; 46(11):1331–1336. [PubMed: 8229111]

4. Maria VA, Victorino RM. Development and validation of a clinical scale for the diagnosis of drug-induced hepatitis. Hepatology. 1997 Sep; 26(3):664–669. [PubMed: 9303497]

5. Rochon J, Protiva P, Seeff LB, Fontana RJ, Liangpunsakul S, Watkins PB, et al. Reliability of the Roussel Uclaf Causality Assessment Method for assessing causality in drug-induced liver injury. Hepatology. 2008 Oct; 48(4):1175–1183. [PubMed: 18798340]

6. Fontana RJ, Watkins PB, Bonkovsky HL, Chalasani N, Davern T, Serrano J, et al. Drug-Induced Liver Injury Network (DILIN) prospective study: rationale, design and conduct. Drug Saf. 2009; 32(1):55–68. [PubMed: 19132805]

7. Agarwal VK, McHutchison JG, Hoofnagle JH. Important elements for the diagnosis of drug-induced liver injury. Clin Gastroenterol Hepatol. 2010 May; 8(5):463–470. [PubMed: 20170750]

8. Davern TJ, Chalasani N, Fontana RJ, Hayashi PH, Protiva P, Kleiner DE, et al. Acute hepatitis E infection accounts for some cases of suspected drug-induced liver injury. Gastroenterology. 2011 Nov; 141(5):1665–1672. e1–e9. [PubMed: 21855518]

9. Garcia-Cortes M, Stephens C, Lucena MI, Fernandez-Castaner A, Andrade RJ. Causality assessment methods in drug induced liver injury: strengths and weaknesses. J Hepatol. 2011 Sep; 55(3):683–691. [PubMed: 21349301]

10. Rousselet MC, Michalak S, Dupre F, Croue A, Bedossa P, Saint-Andre JP, et al. Sources of variability in histological scoring of chronic viral hepatitis. Hepatology. 2005 Feb; 41(2):257–264. [PubMed: 15660389]

11. Kleiner DE, Brunt EM, Van Natta M, Behling C, Contos MJ, Cummings OW, et al. Design and validation of a histological scoring system for nonalcoholic fatty liver disease. Hepatology. 2005 Jun; 41(6):1313–1321. [PubMed: 15915461]

12. Dalton HR, Fellows HJ, Stableforth W, Joseph M, Thurairajah PH, Warshow U, et al. The role of hepatitis E virus testing in drug-induced liver injury. Aliment Pharmacol Ther. 2007 Nov 15; 26(10):1429–1435. [PubMed: 17850420]

13. Chalasani N, Fontana RJ, Bonkovsky HL, Watkins PB, Davern T, Serrano J, et al. Causes, clinical features, and outcomes from a prospective study of drug-induced liver injury in the United States. Gastroenterology. 2008 Dec; 135(6):1924–1934. 34 e1–34 e4. [PubMed: 18955056]

14. Zimmerman, HJ. Antituberculosis agents. In: Zimmerman, HJ., editor. Hepatotoxicity: the adverse effects of drugs and other chemicals on the liver. 2nd ed. Philadelphia: Lipponcott Williams & Wilkins; 1999. p. 611-621.

15. Lucena MI, Andrade RJ, Fernandez MC, Pachkoria K, Pelaez G, Duran JA, et al. Determinants of the clinical expression of amoxicillin-clavulanate hepatotoxicity: a prospective series from Spain. Hepatology. 2006 Oct; 44(4):850–856. [PubMed: 17006920]

16. Verma, S.; Kaplowitz, N. Hepatotoxicity of antituberculosis drugs. In: Kaplowitz, N.; DeLeve, LD., editors. Drug-induced liver disease. 2nd ed.. New York: Informa Healthcare USA; 2007. p. 547-566.

17. Aithal GP, Rawlins MD, Day CP. Accuracy of hepatic adverse drug reaction reporting in one English health region. BMJ. 1999 Dec 11.319(7224):1541. [PubMed: 10591713]

18. Teschke R, Schulze J, Schwarzenboeck A, Eickhoff A, Frenzel C. Herbal hepatotoxicity: suspected cases assessed for alternative causes. Eur J Gastroenterol Hepatol. 2013 Sep; 25(9):1093–1098. [PubMed: 23510966]

19. Hennes EM, Zeniya M, Czaja AJ, Pares A, Dalekos GN, Krawitt EL, et al. Simplified criteria for the diagnosis of autoimmune hepatitis. Hepatology. 2008 Jul; 48(1):169–176. [PubMed: 18537184]

20. Agarwal VK, McHutchison JG, Hoofnagle JH. Important Elements for the Diagnosis of Drug-Induced Liver Injury. Clin Gastroenterol Hepatol. 2010 Feb 17.

21. Teschke R, Eickhoff A, Schulze J. Drug- and herb-induced liver injury in clinical and translational hepatology: causality assessment methods, quo vadis? Journal of Clinical and Translational Hepatology. 2013; 1:59–74.

a:

Group A

Group B

41%

51%

4%

4%

35%

59%

6%

5-point likelihood scoring scale

☐ Same score
☒ 1 pt. difference
◩ 2 pt. difference
■ 3 pt. difference

b:

Group A

Group B

35%

57%

6%

2%

4-point likelihood scoring scale

25%

2%

73%

**Figure 1.**
Magnitude of differences between original and reassessment scores expressed as
percentages of cases with 0, 1, 2 and 3 point differences. There were no cases differing by 4
points. Groups A and B shown separately. (a) Percentages for the 5-point likelihood scoring
scale. (b). Percentages for a 4-point likelihood scoring scale that combines scores 1 and 2.
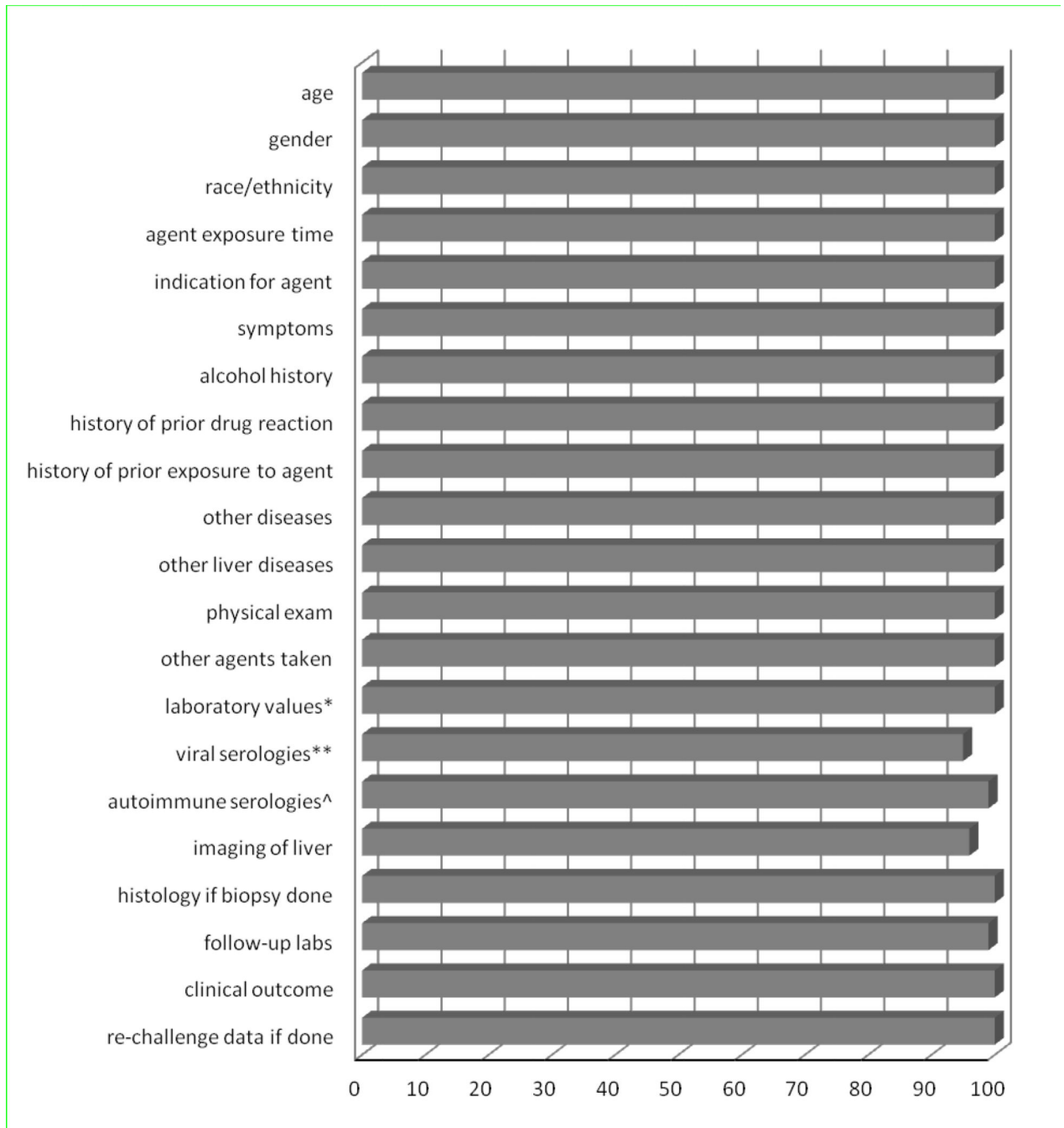
**Figure (Appendix).**

Completeness of data at adjudication for parameters adapted from Agarwal et al.(7)

* laboratory values including liver biochemistries, INR, cell count at onset and enrollment

**complete viral serologies for acute hepatitis A, B and C

^ autoimmune serologies (ANA, ASMA)

**Table 1**

Features of DILIN cases randomly selected for reassessment compared to cases not reassessed.

| Characteristic | Cases Reassessed N=100 | Cases Not Reassessed N=983 | p-value |
|---|---|---|---|
| **Demographics:** | | | |
| Age | | | |
|   Mean (SD) | 50.0 (18.59) | 48.8 (17.01) | 0.41 |
| Gender | | | |
|   Female | 63 (63.0%) | 558 (56.8%) | 0.24 |
| Race (self report) | | | |
|   Caucasian | 78/100 (78.0%) | 764/973 (78.5%) | |
|   Black | 12/100 (12.0%) | 113/973 (11.6%) | 0.96 |
|   Asian | 3/100 (3.0%) | 36/973 (3.7%) | |
|   Other/Multiracial | 7/100 (7.0%) | 60/973 (6.2%) | |
| Body mass index (BMI) | | | |
|   number w/ BMI | 88 | 917 | |
|   Mean (SD) | 27.0 (6.89) | 27.4 (6.50) | 0.39 |
| **Liver Injury** | | | |
| Categorized Days from Primary Drug Start to DILI Onset | | | |
|   <= 1 week | 7/88 (8.0%) | 81/813 (10.0%) | |
|   2 to 4 weeks | 29/88 (33.0%) | 255/813 (31.4%) | |
|   5 to 12 weeks | 26/88 (29.5%) | 272/813 (33.5%) | 0.84 |
|   13 to 24 weeks | 11/88 (12.5%) | 85/813 (10.5%) | |
|   > 24 weeks | 15/88 (17.0%) | 120/813 (14.8%) | |
| # of concomitant drugs in the 2 months prior to DILI onset | | | |
|   0 to 2 | 24/93 (25.8%) | 191/824 (23.2%) | |
|   3 to 5 | 25/93 (26.9%) | 223/824 (27.1%) | 0.83 |
|   > 5 | 44/93 (47.3%) | 410/824 (49.8%) | |
| Pattern of liver injury at onset or earliest after onset[*] | | | |
|   Cholestatic (R < 2) | 25/100 (25.0%) | 248/974 (25.5%) | |
|   Mixed (R 2–5) | 17/100 (17.0%) | 207/974 (21.3%) | 0.59 |
|   Hepatocellular (R > 5) | 58/100 (58.0%) | 519/974 (53.3%) | |
| Peak values between DILI onset and 6 mo. after enrollment | | | |
|   ALT (U/L) | | | |
|     Mean (SD) | 1118 (1887) | 998 (1481) | 0.90 |
|   Alkaline Phosphatase (U/L) | | | |
|     Mean (SD) | 390 (381) | 410 (418) | 0.33 |
|   Total Bilirubin (mg/dL) | | | |
|     Median (25th,75th) | 9.0 (1.7, 22.9) | 9.7 (2.7, 19.5) | 0.94 |
|   INR | | | |
|     Mean (SD) | 2.1 (2.49) | 1.6 (1.41) | 0.07 |

| Characteristic | Cases Reassessed N=100 | Cases Not Reassessed N=983 | p-value |
|---|---|---|---|
| **Adjudication** | | | |
| Overall causality score | | | |
| Definite Greater than 95% | 23/100 (23.0%) | 173/744 (23.3%) | |
| Very likely 75–95% | 39/100 (39.0%) | 312/744 (41.9%) | |
| Probable 50–75% | 19/100 (19.0%) | 141/744 (19.0%) | 0.22 |
| Possible 25–50% | 10/100 (10.0%) | 91/744 (12.2%) | |
| Unlikely Less than 25% | 9/100 (9.0%) | 27/744 (3.6%) | |
| **Outcomes** | | | |
| DILIN severity score(6) | | | |
| Mild | 23/100 (23.0%) | 185/744 (24.9%) | |
| Moderate | 19/100 (19.0%) | 156/744 (21.0%) | |
| Moderate-hospitalized | 35/100 (35.0%) | 221/744 (29.7%) | 0.36 |
| Severe | 12/100 (12.0%) | 129/744 (17.3%) | |
| Fatal | 11/100 (11.0%) | 53/744 (7.1%) | |

[*]
$R = (ALT/ULN) \div (AP/ULN)$

**Table 2**

Original and reassessment score frequencies and kappa statistics for Groups A, B and total cohort. (a) 5-point likelihood scale (b) 4-point likelihood scale, combining categories 1 and 2.

**(a).**

| Original Score Frequencies | | | Reassessment Score Frequencies | | | | | Totals | Kappa (95% CI) |
|---|---|---|---|---|---|---|---|---|---|
| | | | 1 | 2 | 3 | 4 | 5 | | |
| Group A n = 49 | 1 | | 9 | 3 | 0 | 1 | 0 | 13 | 0.53 (0.38–0.68) |
| | 2 | | 0 | 12 | 6 | 1 | 1 | 20 | |
| | 3 | | 0 | 3 | 1 | 3 | 1 | 8 | |
| | 4 | | 0 | 0 | 1 | 3 | 1 | 5 | |
| | 5 | | 0 | 0 | 0 | 3 | 0 | 3 | |
| Group B n= 51 | 1 | | 6 | 3 | 1 | 0 | 0 | 10 | 0.67 (0.53–0.81) |
| | 2 | | 4 | 10 | 4 | 1 | 0 | 19 | |
| | 3 | | 1 | 2 | 4 | 4 | 0 | 11 | |
| | 4 | | 0 | 0 | 1 | 4 | 0 | 5 | |
| | 5 | | 0 | 0 | 0 | 0 | 6 | 6 | |
| Groups A & B n = 100 | 1 | | 15 | 6 | 1 | 1 | 0 | 23 | 0.60 (0.50–0.71) |
| | 2 | | 4 | 22 | 10 | 2 | 1 | 39 | |
| | 3 | | 1 | 5 | 5 | 7 | 1 | 19 | |
| | 4 | | 0 | 0 | 2 | 7 | 1 | 10 | |
| | 5 | | 0 | 0 | 0 | 3 | 6 | 9 | |

**(b)**

| Original Score Frequencies | | Reassessment Score Frequencies | | | | Totals | Kappa (95% CI) |
|---|---|---|---|---|---|---|---|
| | | 1 or 2 | 3 | 4 | 5 | | |
| Group A n = 49 | 1 or 2 | 24 | 6 | 2 | 1 | 33 | 0.44 (0.27–0.61) |
| | 3 | 3 | 1 | 3 | 1 | 8 | |
| | 4 | 0 | 1 | 3 | 1 | 5 | |
| | 5 | 0 | 0 | 3 | 0 | 3 | |
| Group B n = 51 | 1 or 2 | 23 | 5 | 1 | 0 | 29 | 0.73 (0.58–0.87) |
| | 3 | 3 | 4 | 4 | 0 | 11 | |
| | 4 | 0 | 1 | 4 | 0 | 5 | |
| | 5 | 0 | 0 | 0 | 6 | 6 | |
| Groups A & B n = 100 | 1 or 2 | 47 | 11 | 3 | 1 | 62 | 0.6 (0.48–0.72) |
| | 3 | 6 | 5 | 7 | 1 | 19 | |
| | 4 | 0 | 2 | 7 | 1 | 10 | |
| | 5 | 0 | 0 | 3 | 6 | 9 | |

1 = definite (>95% likelihood), 2 = very likely (75–95%), 3 = probable (50–74%), 4 = possible (25–49%), 5 = unlikely (<25%)

**Table 3**

Direction of score changes on reassessment stratified by original score.

| Original Score | N | Less likely DILI | More likely DILI | No change |
|:---:|:---:|:---:|:---:|:---:|
| 1 | 23 | 35% | NA[*] | 65% |
| 2 | 39 | 33% | 10% | 57% |
| 3 | 19 | 42% | 32% | 26% |
| 4 | 10 | 10% | 20% | 70% |
| 5 | 9 | NA[*] | 33% | 67% |

[*] NA = not applicable

**Table 4**

Agents, scores and causes for DILI diagnosis uncertainty in cases crossing the 50% likelihood threshold on reassessment.

| # | Agents | Initial Score | Reassessment Score | Potential Reasons for DILI Diagnosis Uncertainty | |
|---|--------|---------------|--------------------|--------------------------------------------------|--|
| 1 | Ranitidine | 2 | 4 | Timing<br>Competing diagnosis | Short latency (2 days)<br>Choledocholithiasis |
| 2 | Duloxetine | 1 | 4 | Competing diagnosis | Biliary obstruction |
| 3 | Topiramate | 2 | 5 | Competing diagnosis | Hepatitis C, acute |
| 4 | Cefuroxime/Nystatin | 3 | 4 | Competing diagnosis | Autoimmune hepatitis |
| 5 | Cephalexin/Levofloxacin | 3 | 5 | Competing diagnosis | Hepatitis C, chronic |
| 6 | Experimental agent | 4 | 3 | Competing diagnosis | Ischemic hepatopathy |
| 7 | Ezetimibe/Simvastatin | 3 | 4 | Competing diagnosis | Occult alcohol |
| 8 | Sulfamethoxazole-Trimethoprim | 2 | 4 | Timing | Unclear timing from drug exposure to elevated liver enzymes |
| 9 | Ceftriaxone/Ampicillin-Sulbactam/Fluconazole | 3 | 4 | Competing diagnosis | Ischemic hepatopathy; occult acetaminophen overdose |
| 10 | Ciprofloxacin | 4 | 3 | Competing diagnosis | Hepatitis, atypical viral (giant cell) |
| 11 | Hydroxycut$_{TM}$ (HDS*) | 3 | 4 | Timing | Long latency (approximately 75 days) after stopping agent |
| 12 | Azithromycin | 3 | 4 | Competing diagnosis | Ischemic hepatopathy; occult acetaminophen overdose |
| 13 | Ultra Vitality$_{TM}$ (HDS*) | 3 | 4 | Timing | Long latency (approximately 2190 days) while on agent |
| 14 | 6-Mercaptopurine | 3 | 4 | Timing | Unclear timing from drug exposure to elevated liver enzymes |

*
Herbal Dietary Supplement